

Queue management

Outline

◆ Characteristics of a Waiting-Line System.

- ◆ Arrival characteristics.
- ◆ Waiting-Line characteristics.
- ◆ Service facility characteristics.

◆ Waiting Line (Queuing) Models.

- ◆ M/M/1: One server.
- ◆ M/M/2: Two servers.
- ◆ M/M/S: S servers.
- ◆ Cost comparisons.

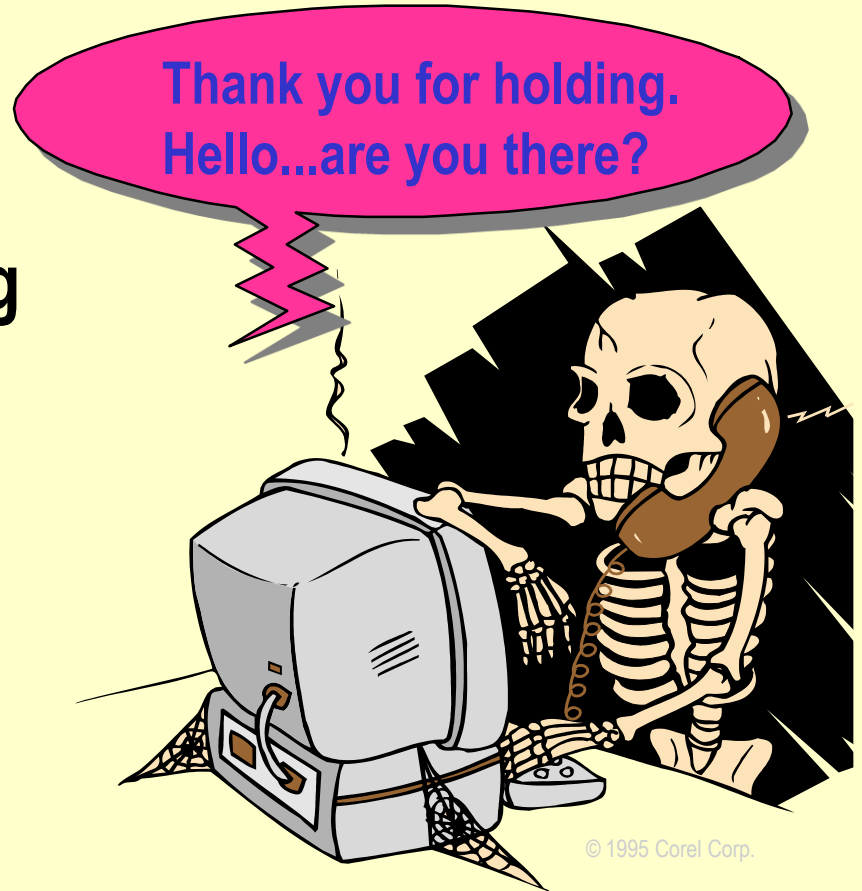
Waiting Lines

- ◆ **First studied by A. K. Erlang in 1913.**
 - ◆ **Analyzed telephone facilities.**
- ◆ **Body of knowledge called queuing theory.**
 - ◆ **Queue is another name for waiting line.**
- ◆ **Decision problem:**
 - ◆ **Balance cost of providing good service with cost of customers waiting.**

You've Been There Before!

The average person
spends 5 years waiting
in line!!

'The other line
always moves faster.'



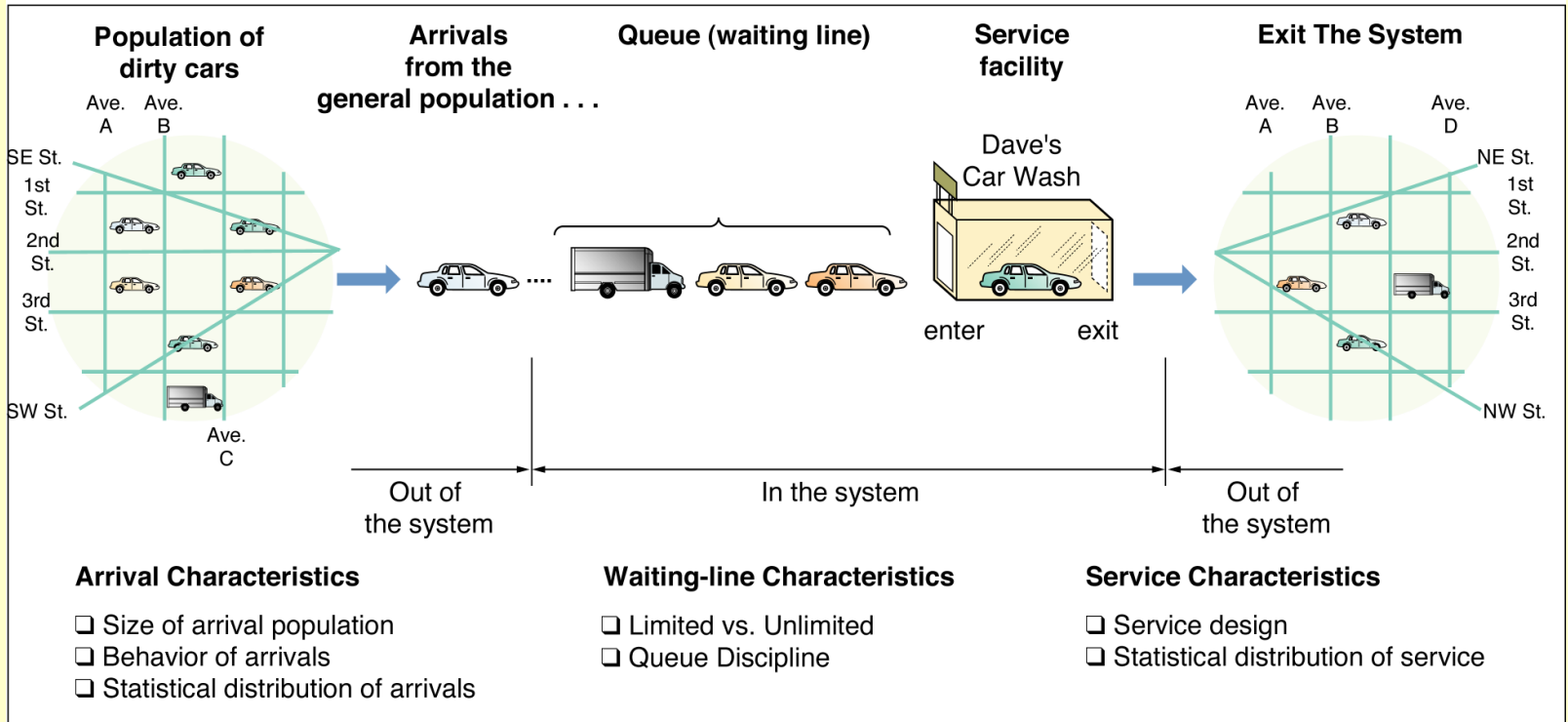
Waiting Line Examples

<u>Situation</u>	<u>Arrivals</u>	<u>Servers</u>	<u>Service Process</u>
Bank	Customers	Teller	Deposit etc.
Doctor's office	Patient	Doctor	Treatment
Traffic intersection	Cars	Traffic Signal	Controlled passage
Assembly line	Parts	Workers	Assembly

Waiting Line Components

- ◆ **Arrivals:** Customers (people, machines, calls, etc.) that demand service.
- ◆ **Waiting Line (Queue):** Arrivals waiting for a free server.
- ◆ **Servers:** People or machines that provide service to the arrivals.
- ◆ **Service System:** Includes waiting line and servers.

Car Wash Example



Key Tradeoff

- ◆ Higher service level (more servers, faster servers)
 - ⇒ Higher costs to provide service.
 - ⇒ Lower cost for customers waiting in line (less waiting time).

Waiting Line Terminology

- ◆ **Queue:** Waiting line.
- ◆ **Arrival:** 1 person, machine, part, etc. that arrives and demands service.
- ◆ **Queue discipline:** Rules for determining the order that arrivals receive service.
- ◆ **Channels:** Parallel servers.
- ◆ **Phases:** Sequential stages in service.

Input Characteristics

◆ Input source (population) size.

- ◆ **Infinite**: Number in service does not affect probability of a new arrival.
 - ◆ A very large population can be treated as infinite.
- ◆ **Finite**: Number in service affects probability of a new arrival.
 - ◆ Example: Population = 10 aircraft that may need repair.

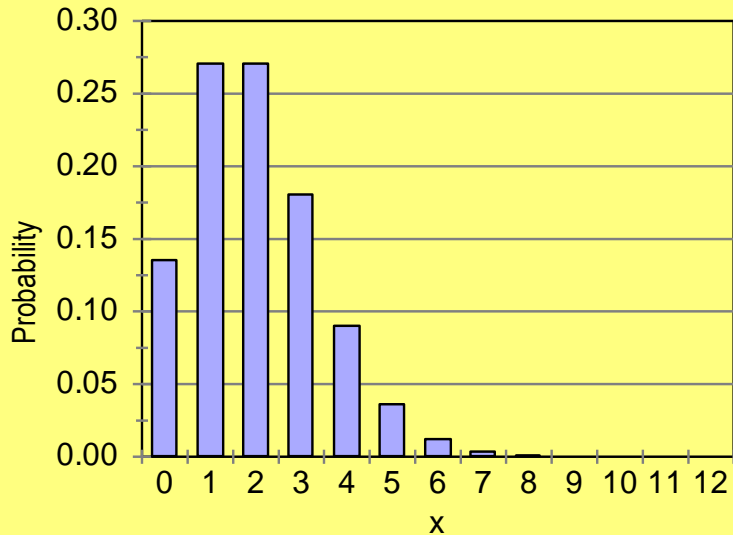
◆ Arrival pattern.

- ◆ **Random**: Use Poisson probability distribution.
- ◆ **Non-random**: Appointments.

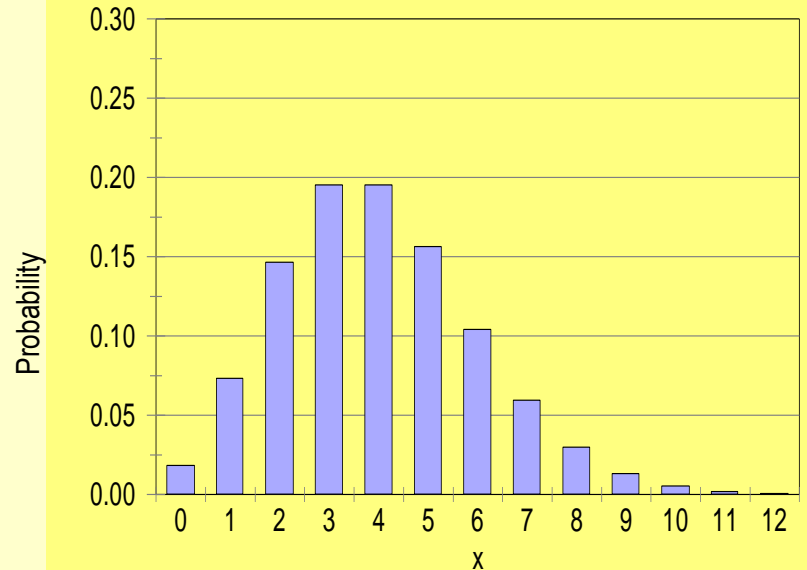
Poisson Distribution

- ◆ Number of events that occur in an interval of time.
 - ◆ Example: Number of customers that arrive each half-hour.
- ◆ Discrete distribution with mean = λ
 - ◆ Example: Mean arrival rate = 5/hour .
 - ◆ Probability:
$$P(x) = \frac{e^{-\lambda} \lambda^x}{x!}$$
- ◆ Time between arrivals has a negative exponential distribution.

Poisson Probability Distribution



$\lambda=2$



$\lambda=4$

Behavior of Arrivals

- ◆ **Patient.**

- ◆ Arrivals will wait in line for service.

- ◆ **Impatient.**

- ◆ **Balk:** Arrival leaves before entering line.
 - ◆ Arrival sees long line and decides to leave.
 - ◆ **Renegue:** Arrival leaves after waiting in line a while.

Waiting Line Characteristics

◆ Line length:

- ◆ **Limited:** Maximum number waiting is limited.
 - ◆ Example: Limited space for waiting.
- ◆ **Unlimited:** No limit on number waiting.

◆ Queue discipline:

- ◆ **FIFO (FCFS):** First in, First out. (First come, first served).
- ◆ **Random:** Select arrival to serve at random from those waiting.
- ◆ **Priority:** Give some arrivals priority for service.

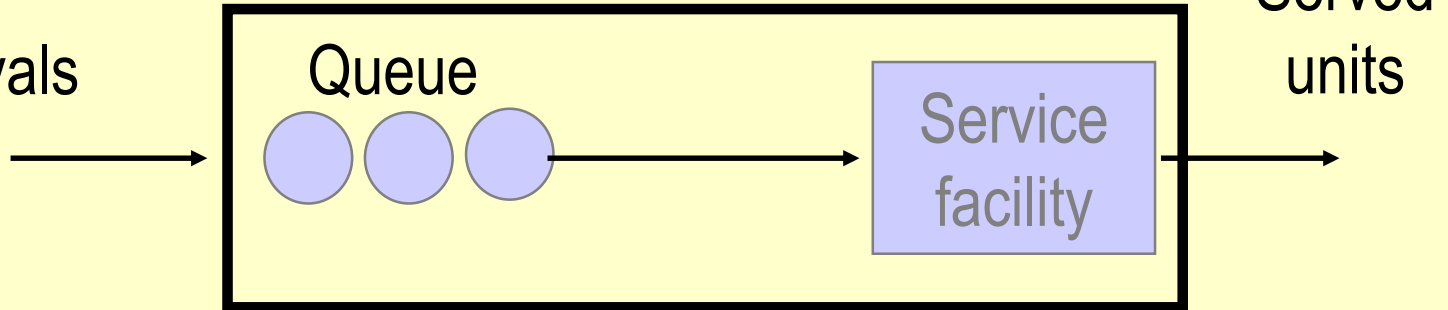
Service Configuration

- ◆ **Single channel, single phase.**
 - ◆ One server, one phase of service.
- ◆ **Single channel, multi-phase.**
 - ◆ One server, multiple phases in service.
- ◆ **Multi-channel, single phase.**
 - ◆ Multiple servers, one phase of service.
- ◆ **Multi-channel, multi-phase.**
 - ◆ Multiple servers, multiple phases of service.

Single Channel, Single Phase

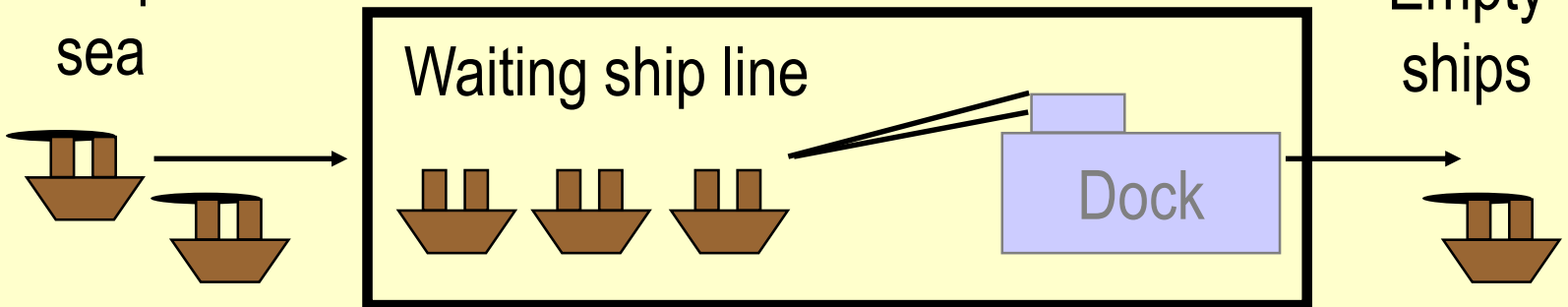
Service system

Arrivals



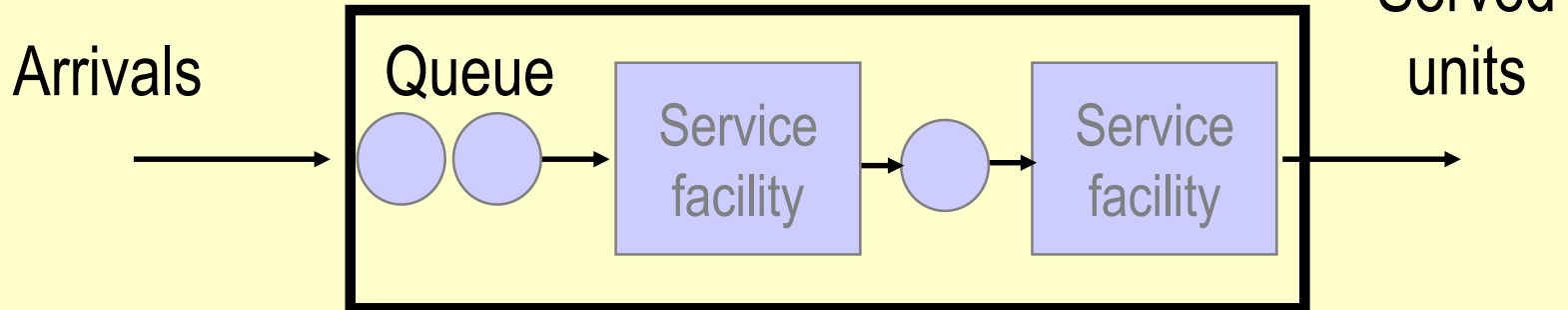
Ships at sea

Ship unloading system

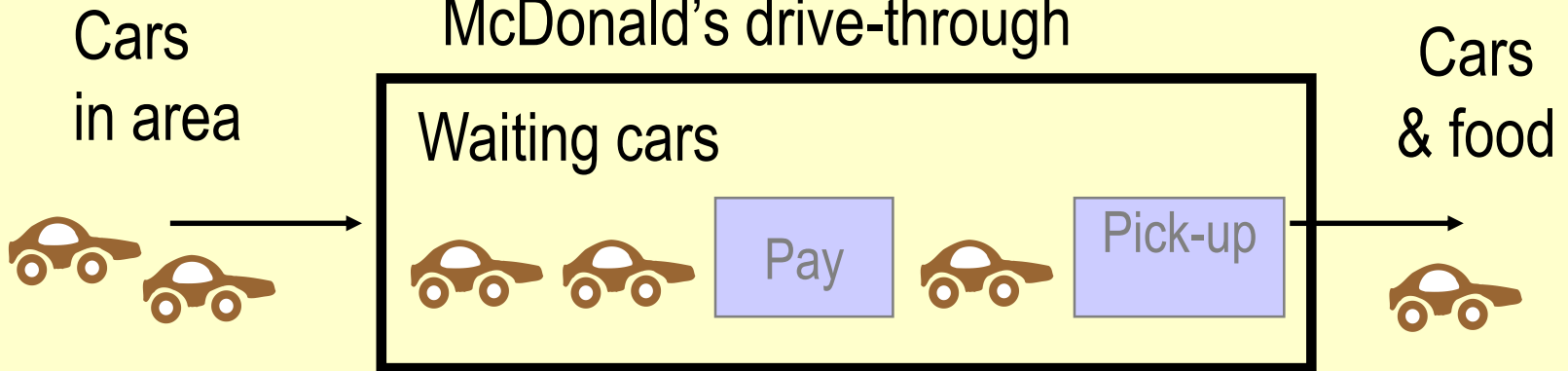


Single Channel, Multi-Phase

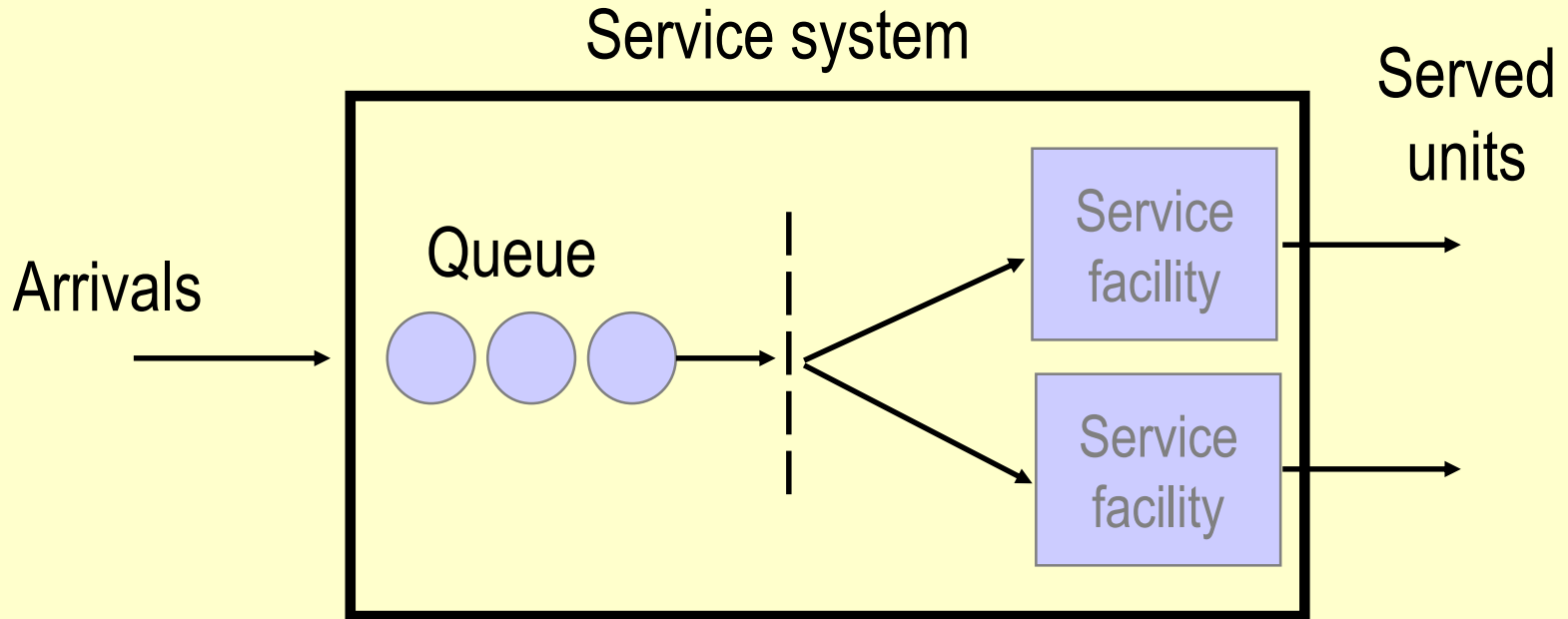
Service system



McDonald's drive-through

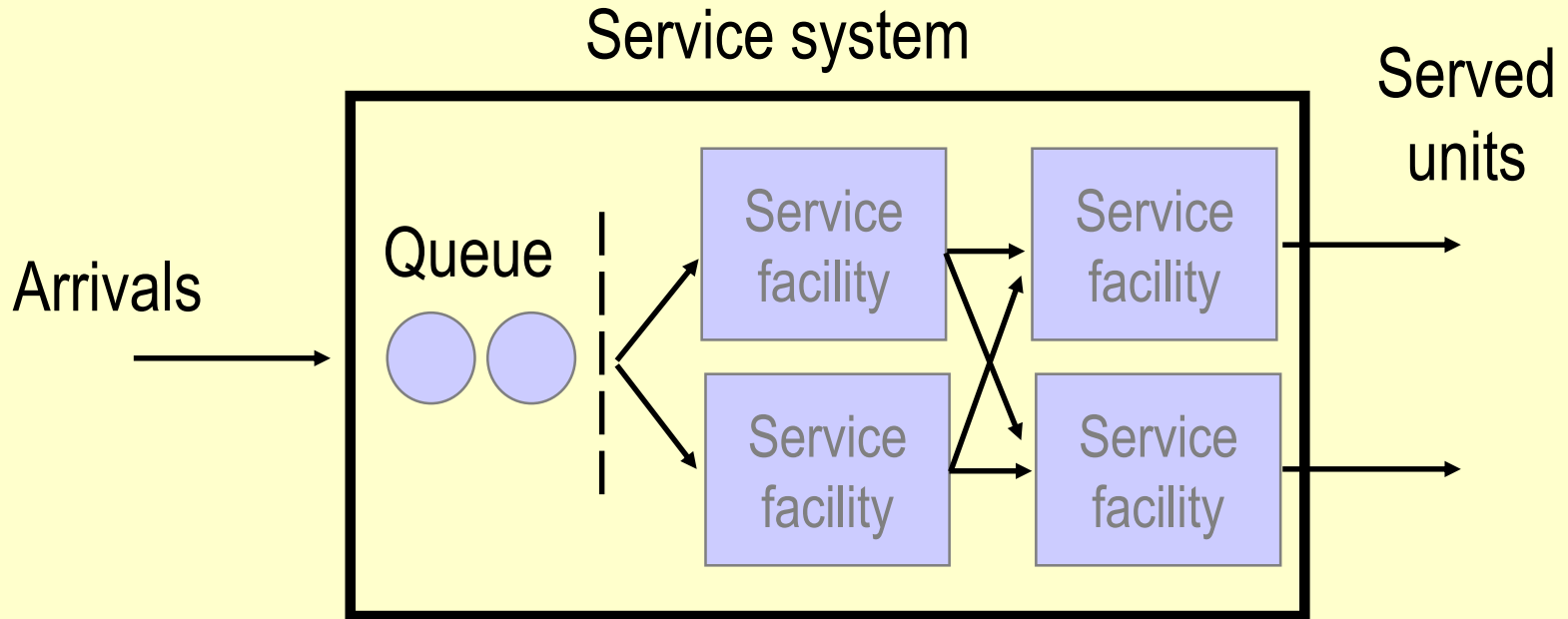


Multi-Channel, Single Phase



Example: Bank customers wait in ***single line*** for one of several tellers.

Multi-Channel, Multi-Phase



Example: At a laundromat, customers use one of several washers, then one of several dryers.

Service Times

- ◆ **Random:** Use Negative exponential probability distribution.
 - ◆ Mean service *rate* = μ
 - ◆ 6 customers/hr.
 - ◆ Mean service *time* = $1/\mu$
 - ◆ 1/6 hour = 10 minutes.
- ◆ **Non-random:** May be constant.
 - ◆ Example: Automated car wash.

Negative Exponential Distribution

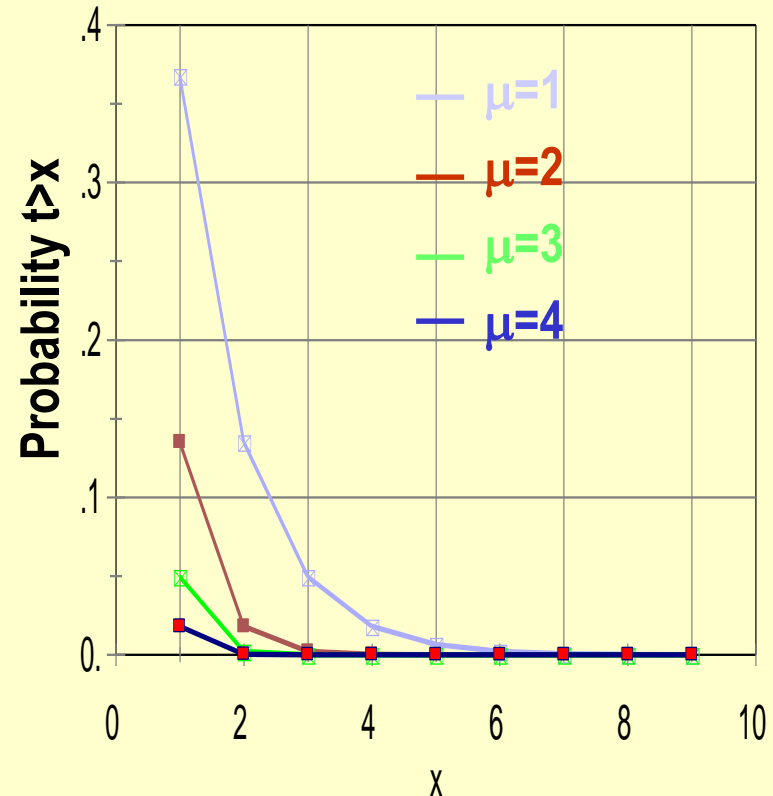
- ◆ Continuous distribution.

- ◆ Probability:

$$f(t > x) = e^{-\mu x}$$

- ◆ Example: Time between arrivals.

- ◆ Mean service **rate** = μ
 - ◆ 6 customers/hr.
 - ◆ Mean service **time** = $1/\mu$
 - ◆ 1/6 hour = 10 minutes



Assumptions in the Basic Model

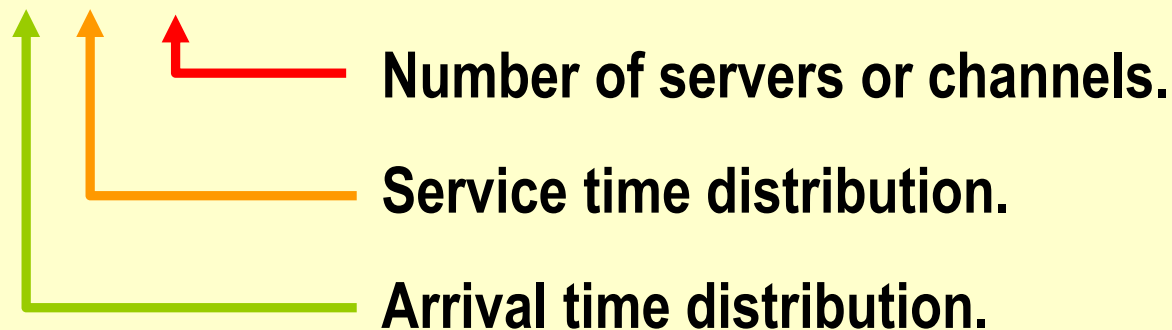
- ◆ Customer population is *homogeneous* and *infinite*.
- ◆ Queue capacity is *infinite*.
- ◆ Customers are *well behaved* (no balking or reneging).
- ◆ Arrivals are served *FCFS (FIFO)*.
- ◆ *Poisson arrivals*.
 - ◆ The time between arrivals follows a negative exponential distribution
- ◆ *Exponential service times*: Services are described by the negative exponential distribution.

Steady State Assumptions

- ◆ Mean arrival rate λ , mean service rate μ , and the number of servers are constant.
- ◆ The service rate is greater than the arrival rate.
- ◆ These conditions have existed for a long time.

Queuing Model Notation

a/b/S



- ◆ M = Negative exponential distribution (Poisson arrivals).
- ◆ G = General distribution.
- ◆ D = Deterministic (scheduled).

Types of Queuing Models

- ◆ **Simple (M/M/1).**

- ◆ **Example: Information booth at mall.**

- ◆ **Multi-channel (M/M/S).**

- ◆ **Example: Airline ticket counter.**

- ◆ **Constant Service (M/D/1).**

- ◆ **Example: Automated car wash.**

- ◆ **Limited Population.**

- ◆ **Example: Department with only 7 drills that may break down and require service.**

Common Questions

- ◆ Given λ , μ and S , how large is the queue (waiting line)?
- ◆ Given λ and μ , how many servers (channels) are needed to keep the average wait within certain limits?
- ◆ What is the total cost for servers and customer waiting time?
- ◆ Given λ and μ , how many servers (channels) are needed to minimize the total cost?

Performance Measures

- ◆ Average queue time = W_q
- ◆ Average queue length = L_q
- ◆ Average time in system = W_s
- ◆ Average number in system = L_s
- ◆ Probability of idle service facility = P_0
- ◆ System utilization = ρ
- ◆ Probability of more than k units in system = $P_{n > k}$
 - ◆ Also, fraction of time there are more than k units in the system.

General Queuing Equations

$$\rho = \frac{\lambda}{S\mu}$$

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

$$L_s = L_q + \frac{\lambda}{\mu}$$

$$L_q = \lambda W_q$$

$$L_s = \lambda W_s$$

Given one of W_s , W_q , L_s , or L_q you can use these equations to find all the others.

M/M/1 Model

- ◆ **Type:** Single server, single phase system.
- ◆ **Input source:** Infinite; no balks, no reneging.
- ◆ **Queue:** Unlimited; single line; FIFO (FCFS).
- ◆ **Arrival distribution:** Poisson.
- ◆ **Service distribution:** Negative exponential.

M/M/1 Model Equations

Average # of customers in system: $L_s = \frac{\lambda}{\mu - \lambda}$

Average time in system: $W_s = \frac{1}{\mu - \lambda}$

Average # of customers in queue: $L_q = \frac{\lambda^2}{\mu(\mu - \lambda)}$

Average time in queue: $W_q = \frac{\lambda}{\mu(\mu - \lambda)}$

System utilization $\rho = \frac{\lambda}{\mu}$

M/M/1 Probability Equations

Probability of 0 units in system, i.e., system idle:

$$P_0 = 1 - \rho = 1 - \frac{\lambda}{\mu}$$

Probability of more than k units in system:

$$P_{n>k} = \left(\frac{\lambda}{\mu}\right)^{k+1}$$

This is also probability of k+1 or more units in system.

M/M/1 Example 1

Average arrival rate is 10 per hour. Average service time is 5 minutes.

$$\lambda = 10/\text{hr} \quad \text{and} \quad \mu = 12/\text{hr}$$
$$(1/\mu = 5 \text{ minutes} = 1/12 \text{ hour})$$

Q1: What is the average time between departures?
5 minutes? 6 minutes?

Q2: What is the average wait in the system?

$$W_s = \frac{1}{12/\text{hr} - 10/\text{hr}} = 0.5 \text{ hour or } 30 \text{ minutes}$$

M/M/1 Example 1

$\lambda = 10/\text{hr}$ and $\mu = 12/\text{hr}$

Q3: What is the average wait in line?

$$W_q = \frac{10}{12(12-10)} = 0.41667 \text{ hours} = 25 \text{ minutes}$$

Also note: $W_s = W_q + \frac{1}{\mu}$

so $W_q = W_s - \frac{1}{\mu} = \frac{1}{2} - \frac{1}{12} = 0.41667 \text{ hours}$

M/M/1 Example 1

$\lambda = 10/\text{hr}$ and $\mu = 12/\text{hr}$

Q4: What is the average number of customers in line and in the system?

$$L_q = \frac{10^2}{12(12-10)} = 4.1667 \text{ customers}$$

$$L_s = \frac{10}{12-10} = 5 \text{ customers}$$

Also note: $L_q = \lambda W_q = 10 \times 0.41667 = 4.1667$

$$L_s = \lambda W_s = 10 \times 0.5 = 5$$

M/M/1 Example 1

$\lambda = 10/\text{hr}$ and $\mu = 12/\text{hr}$

Q5: What is the fraction of time the system is empty (server is idle)?

$$P_0 = 1 - \rho = 1 - \frac{\lambda}{\mu} = 1 - \frac{10}{12} = 16.67\% \text{ of the time}$$

Q6: What is the fraction of time there are more than 5 customers in the system?

$$P_{n>5} = \left(\frac{10}{12}\right)^6 = 33.5\% \text{ of the time}$$

More than 5 in the system...

Note that “more than 5 in the system” is the same as:

- ◆ “more than 4 in line”
- ◆ “5 or more in line”
- ◆ “6 or more in the system”.

All are $P_{n>5}$

M/M/1 Example 1

$\lambda = 10/\text{hr}$ and $\mu = 12/\text{hr}$

Q7: How much time per day (8 hours) are there 5 or more customers in line?

$P_{n>5} = 0.335$ so 33.5% of time there are 6 or more in line.

$0.335 \times 480 \text{ min./day} = 160.8 \text{ min.} = \sim 2 \text{ hr } 40 \text{ min.}$

Q8: What fraction of time are there 3 or fewer customers in line?

$$1 - P_{n>4} = 1 - \left(\frac{10}{12}\right)^5 = 1 - 0.402 = 0.598 \text{ or } 59.8\%$$

M/M/1 Example 2

Five copy machines break down at UM St. Louis per eight hour day on average. The average service time for repair is one hour and 15 minutes.

$$\lambda = 5/\text{day} \quad (\lambda = 0.625/\text{hour})$$

$$1/\mu = 1.25 \text{ hours} = 0.15625 \text{ days}$$

$$\mu = 1 \text{ every } 1.25 \text{ hours} = 6.4/\text{day}$$

Q1: What is the average number of “customers” in the system?

$$L_s = \frac{5/\text{day}}{6.4/\text{day} - 5/\text{day}} = 3.57 \text{ broken copiers}$$

M/M/1 Example 2

$$\lambda = 5/\text{day} \quad (\text{or } \lambda = 0.625/\text{hour})$$

$$\mu = 6.4/\text{day} \quad (\text{or } \mu = 0.8/\text{hour})$$

Q2: How long is the average wait in line?

$$W_q = \frac{5}{6.4(6.4 - 5)} = 0.558 \text{ days (or 4.46 hours)}$$

$$W_q = \frac{0.625}{0.8(0.8 - 0.625)} = 4.46 \text{ hours}$$

M/M/1 Example 2

$$\lambda = 5/\text{day} \quad (\text{or } \lambda = 0.625/\text{hour})$$

$$\mu = 6.4/\text{day} \quad (\text{or } \mu = 0.8/\text{hour})$$

Q3: How much time per day (on average) are there 2 or more broken copiers waiting for the repair person?

2 or more “in line” = more than 2 in the system

$$P_{n>2} = \left(\frac{5}{6.4} \right)^3 = 0.477 \quad (47.7\% \text{ of the time})$$

$$0.477 \times 480 \text{ min./day} = 229 \text{ min.} = 3 \text{ hr } 49 \text{ min.}$$

M/M/1 Example 3

A coffee shop sees on average one arrival every two minutes in the morning. The average service time (for preparing the drink, paying, etc.) is 90 seconds. If you leave your house at 9:00 am, and it is a 10 minutes drive to the coffee shop, and then a 15 minute drive to school, what time would you expect to arrive at school (on average)?

$$\lambda = 30/\text{hr} \quad \mu = 40/\text{hr} = 1 \text{ every } 90 \text{ seconds}$$

$$W_q = \frac{30}{40(40 - 30)} = 0.075 \text{ hrs} = 4.5 \text{ minutes}$$

Arrival at school on average is at 9:31 am
(9:00 am + 10 min + 4.5 min + 1.5 min + 15 min = 9:31 am)

M/M/S Model

- ◆ **Type:** Multiple servers; single-phase.
- ◆ **Input source:** Infinite; no balks, no reneging.
- ◆ **Queue:** Unlimited; multiple lines; FIFO (FCFS).
- ◆ **Arrival distribution:** Poisson.
- ◆ **Service distribution:** Negative exponential.

M/M/S Equations

Probability of zero people or units in the system:

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

Average number of people or units in the system:

$$L_s = \frac{\lambda \mu \left(\frac{\lambda}{\mu} \right)^M}{(M-1)!(M\mu - \lambda)} P_0 \frac{\lambda}{\mu}$$

Average time a unit spends in the system:

$$W_s = \frac{\mu \left(\frac{\lambda}{\mu} \right)^M}{(M-1)!(M\mu - \lambda)} P_0 + \frac{1}{\mu}$$

Note: M = number of servers in these equations

M/M/S Equations

Average number of people or
units waiting for service:

$$L_q = L_s - \frac{\lambda}{\mu}$$

Average time a person or
unit spends in the queue:

$$W_q = W_s - \frac{1}{\mu}$$

M/M/2 Model Equations

Average time in system: $W_s = \frac{4\mu}{4\mu^2 - \lambda^2}$

Average time in queue: $W_q = \frac{\lambda^2}{\mu(2\mu + \lambda)(2\mu - \lambda)}$

Average # of customers in queue: $L_q = \lambda W_q$

Average # of customers in system: $L_s = \lambda W_s$

Probability the system is empty: $P_0 = \frac{2\mu - \lambda}{2\mu + \lambda}$

M/M/2 Example

Average arrival rate is 10 per hour.

Average service time is 5 minutes for each of 2 servers.

$$\lambda = 10/\text{hr}, \mu = 12/\text{hr}, \text{ and } S=2$$

Q1: What is the average wait in the system?

$$W_s = \frac{4 \times 12}{4(12)^2 - (10)^2} = 0.1008 \text{ hours} = 6.05 \text{ minutes}$$

M/M/2 Example

$\lambda = 10/\text{hr}$, $\mu = 12/\text{hr}$, and $S=2$

Q2: What is the average wait in line?

$$W_q = \frac{(10)^2}{12 (2 \times 12 + 10)(2 \times 12 - 10)} = 0.0175 \text{ hrs} = 1.05 \text{ minutes}$$

Also note: $W_s = W_q + \frac{1}{\mu}$

so $W_q = W_s - \frac{1}{\mu} = 0.1008 \text{ hrs} - 0.0833 \text{ hrs} = 0.0175 \text{ hrs}$

M/M/2 Example

$\lambda = 10/\text{hr}$, $\mu = 12/\text{hr}$, and $S=2$

Q3: What is the average number of customers in line and in the system?

$$L_q = \lambda W_q = 10/\text{hr} \times 0.0175 \text{ hr} = 0.175 \text{ customers}$$

$$L_s = \lambda W_s = 10/\text{hr} \times 0.1008 \text{ hr} = 1.008 \text{ customers}$$

M/M/2 Example

$\lambda = 10/\text{hr}$ and $\mu = 12/\text{hr}$

Q4: What is the fraction of time the system is empty (server are idle)?

$$P_0 = \frac{2 \times 12 - 10}{2 \times 12 + 10} = 41.2\% \text{ of the time}$$

M/M/1, M/M/2 and M/M/3

	<u>1 server</u>	<u>2 servers</u>	<u>3 servers</u>
W_q	<i>25 min.</i> <i>= 0.417 hr</i>	<i>1.05 min.</i> <i>= 0.0175 hr</i>	<i>0.1333 min. (8 sec.)</i> <i>= 0.00222 hr</i>
W_s	<i>30 min.</i>	<i>6.05 min.</i>	<i>5.1333 min.</i>
L_q	<i>4.167 cust.</i>	<i>0.175 cust.</i>	<i>0.0222 cust.</i>
L_s	<i>5 cust.</i>	<i>1.01 cust.</i>	<i>0.855 cust.</i>
P_0	<i>16.7%</i>	<i>41.2%</i>	<i>43.2%</i>

Waiting Line Costs

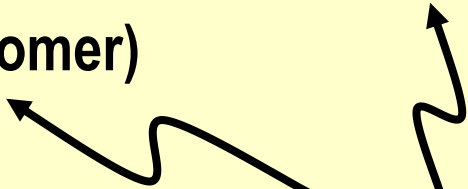
Service cost per day

$$= (\text{\# of servers}) \times (\text{cost per day of each server}) + \\ (\text{\# customers per day}) \times (\text{marginal cost per customer})$$

Customer waiting cost per day

$$= (\text{\# of customers per day}) \times (\text{average wait per customer}) \\ \times (\text{time value for customer})$$

Time units must agree



Service Cost per Day

$\lambda = 10/\text{hr}$ and $\mu = 12/\text{hr}$

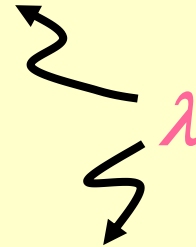
Suppose servers are paid \$7/hr and work 8 hours/day. Also, suppose the marginal cost to serve each customer is \$0.50.

M/M/1 Service cost per day

$$\begin{aligned} &= \$7/\text{hr} \times 8 \text{ hr/day} + \$0.5/\text{cust} \times 10 \text{ cust/hr} \times 8 \text{ hr/day} \\ &= \$56 + \$40 = \$96/\text{day} \end{aligned}$$

M/M/2 Service cost per day

$$\begin{aligned} &= 2 \times \$7/\text{hr} \times 8 \text{ hr/day} + \$0.5/\text{cust} \times 10 \text{ cust/hr} \times 8 \text{ hr/day} \\ &= \$112 + \$40 = \$152/\text{day} \end{aligned}$$



λ

Customer Waiting Cost per Day

$\lambda = 10/\text{hr}$ and $\mu = 12/\text{hr}$

Suppose customer waiting cost is \$10/hr.

M/M/1 Waiting cost per day

$= 0.417 \text{ hr/cust} \times 10 \text{ cust/hr} \times 8 \text{ hr/day} \times \$10/\text{hr} = \$333.33/\text{day}$

M/M/1 total cost per day = \$96 + \$333.33 = \$429.33/day

M/M/2 Waiting cost per day

$= 0.0175 \text{ hr/cust} \times 10 \text{ cust/hr} \times 8 \text{ hr/day} \times \$10/\text{hr} = \$14/\text{day}$

M/M/2 total cost per day = \$152 + \$14 = \$166/day

Unknown Waiting Cost

Suppose customer waiting cost is not known = C .

M/M/1 Waiting cost per day

$$= 0.417 \text{ hr/cust} \times 10 \text{ cust/hr} \times 8 \text{ hr/day} \times C = 33.33C \text{ \$/day}$$

$$\text{M/M/1 total cost per day} = 96 + 33.33C$$

M/M/2 Waiting cost per day

$$= 0.0175 \text{ hr/cust} \times 10 \text{ cust/hr} \times 8 \text{ hr/day} \times C = 1.4C \text{ \$/day}$$

$$\text{M/M/2 total cost per day} = 152 + 1.4C$$

M/M/2 is preferred when $152 + 1.4C < 96 + 33.33C$ or
 $C > \$1.754/\text{hr}$

M/M/2 and M/M/3

Q: How large must customer waiting cost be for M/M/3 to be preferred over M/M/2?

$$M/M/2 \text{ total cost} = 152 + 1.4C$$

M/M/3 Waiting cost per day

$$= C \times 0.00222 \text{ hr/cust} \times 10 \text{ cust/hr} \times 8 \text{ hr/day} = 0.1776C \text{ \$/day}$$

$$M/M/3 \text{ total cost} = 208 + 0.1776C$$

M/M/3 is preferred over M/M/2 when

$$208 + 0.1776C < 152 + 1.4C$$

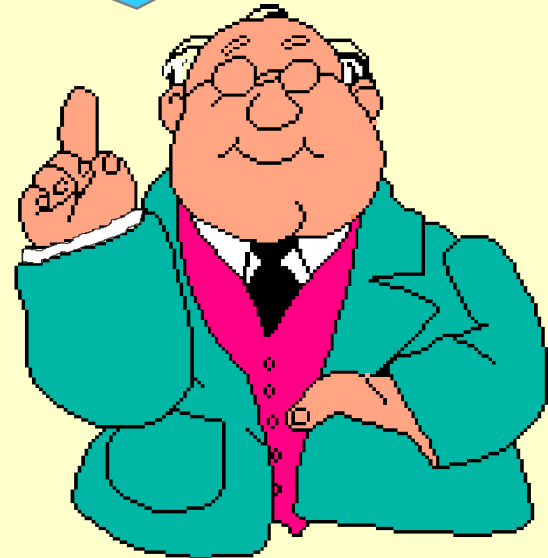
$$C > \$45.81/\text{hr}$$

Remember: λ & μ Are Rates

- ◆ λ = Mean number of arrivals per time period.
 - ◆ Example: 3 units/hour.

If average service time is 15 minutes, then μ is 4 customers/hour

- ◆ μ = Mean number of arrivals served per time period.
 - ◆ Example: 4 units/hour.
 - ◆ $1/\mu = 15$ minutes/unit.



Other Queuing Models

◆ M/D/S

- ◆ Constant service time; Every service time is the same.
- ◆ Random (Poisson) arrivals.

◆ Limited population.

- ◆ Probability of arrival depends on number in service.

◆ Limited queue length.

- ◆ Limited space for waiting.

◆ Many others...

Other Considerations

- ◆ Wait time & queue length increase rapidly for $\rho > 0.7$
 - ◆ Queue is small until system is about 70% busy; then queue grows very quickly.
- ◆ Pooling servers is usually advantageous.
 - ◆ Airport check-in vs. Grocery stores.
- ◆ Variance in arrivals & service times causes long waits.
 - ◆ Long service times cause big waits.
- ◆ Cost of waiting is nonlinear.
 - ◆ Twice as long wait may be more than twice as bad.

More Considerations

- ◆ **Reduce effect of waiting.**
 - ◆ Distract customers with something to do, look at or listen to.
 - ◆ Music, art, mirrors, etc.
 - ◆ Provide feedback on expected length of wait.
 - ◆ “Your call will be answered in 6 minutes...”
- ◆ **Use self service can reduce load on servers.**
 - ◆ Self-service at grocery stores.