

King Fahd University of Petroleum & Minerals Computer Engineering Dept

COE 402 – Computer Systems
Performance Evaluation

Term 043

Dr. Ashraf S. Hasan Mahmoud

Rm 22-148-3

Ext. 1724

Email: ashraf@ccse.kfupm.edu.sa

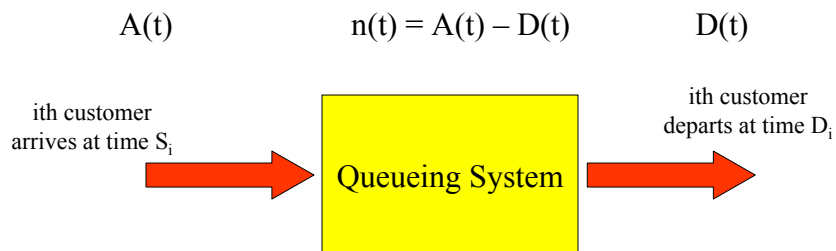
8/20/2005

Dr. Ashraf S. Hasan Mahmoud

1

Queuing Model

- **Consider the following system:**



$$r_i = D_i - A_i$$

$$w_i = r_i - S_i = D_i - A_i - S_i$$

$A(t)$ – number of arrivals in $(0, t]$

$D(t)$ – number of departures in $(0, t]$

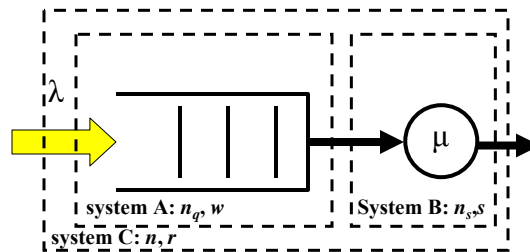
$n(t)$ – number of customers in system in $(0, t]$

r_i – duration of time spent in system for i^{th} customer – **textbook call this response time (r)**

w_i – duration of time spent waiting for service for i^{th} customer

Queuing Model

- The model is characterized using the following quantities:
 - λ = mean arrival rate of customers = $1/E[\tau]$ (remember τ is the interarrival time)
 - s = service time per job
 - μ = mean service rate = $1/E[s]$
 - n = number of job in system $\rightarrow E[n]$
 - n_q = number of jobs in buffer $\rightarrow E[n_q]$
 - ns = number of jobs in server $\rightarrow E[ns]$
 - r = response time or total time in system $\rightarrow E[r]$
 - w = waiting time $\rightarrow E[w]$



8/20/2005

Dr. Ashraf S. Hasan Mahmoud

3

Example: Queueing System

Problem: A data communication line delivers a block of information every 10 microseconds. A decoder check each block for errors and corrects the errors if necessary. It takes 1 microsecond to determine whether the block has any errors. If the block has one error it takes 5 microseconds to correct it and it has more than 1 error it takes 20 microseconds to correct the error. Blocks wait in the queue when the decoder falls behind. Suppose that the decoder is initially empty and that the number of errors in the first 10 blocks are: 0, 1, 3, 1, 0, 4, 0, 1, 0, 0.

- Plot the number of blocks in the decoder as a function of time.
- Find the mean number of blocks in the decoder
- What percent of the time is the decoder empty?

8/20/2005

Dr. Ashraf S. Hasan Mahmoud

4

Example: Queueing System – cont'd

Solution:

Interarrival time = 10 μ sec

Service time = 1 if no errors

1+5 if 1 error

1+20 if more than 1 error

The queue parameters (A, D, S, and W) are shown below:

Block #:	1	2	3	4	5	6	7	8	9	10
Arrivals:	10	20	30	40	50	60	70	80	90	100
Errors:	0	1	3	1	0	4	0	1	0	0
Service:	1	6	21	6	1	21	1	6	1	1
Departs:	11	26	51	57	58	81	82	88	91	101
Waiting:	0	0	0	11	7	0	11	2	0	0

8/20/2005

Dr. Ashraf S. Hasan Mahmoud

5

Example: Queueing System – cont'd

Solution:

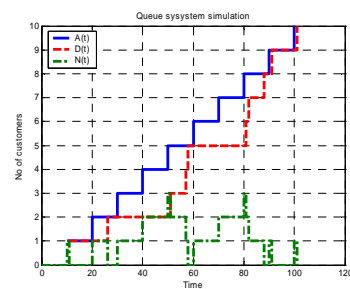
Using the previous results and knowing that

$$n(t) = A(t) - D(t)$$

One can produce the following results

Average no of customers in system = 0.950
 Average customer waiting time = 3.100 microsec
 Maximum simulation time = 101.000 microsec
 Duration server busy = 65.000
 Server utilization = 0.6436
 Server idle = 0.3564

The following Matlab code can be used to solve this queue system (Note the code is general – it solves any system provided The Arrivals vector A, and the service vector S)



8/20/2005

Dr. Ashraf S. Hasan M

Example: Queueing System – cont'd

```

0001 %
0002 % Problem 9.3 - Leon Garcia's book
0003 clear all
0004 A = [10;10;100];
0005 Errors = [0 1 3 1 0 4 0 1 0 0];
0006 S = zeros(size(A));
0007 D = zeros(size(A));
0008 %
0009 % this loop to computes service times
0010 for i=1:length(A)
0011     if (Errors(i)==0) S(i) = 1;
0012     else
0013         if (Errors(i)==1) S(i) = 6;
0014         else
0015             S(i) = 21;
0016         end
0017     end
0018 %
0019 % this section computes the departure time for
the ith user
0020     if (i>1) % this is not the first user
0021         if (D(i-1) < A(i)) D(i) = A(i) + S(i);
0022         else
0023             D(i) = D(i-1) + S(i);
0024         end
0025     else
0026         D(i) = A(i)+S(i);
0027     end
0028 %
0029 % compute waiting time
0030     W(i) = D(i) - A(i) - S(i);
0031 end
0032 %

0033 % Compute N(t)
0034 T = []; % time axis
0035 T(1) = 0; % time origin
0036 N = []; % number of customers
0037 N(1) = 0; % initial condition
0038 k = 2; % place for next insert
0039 A_max = A(length(A)); % last arrival instant
0040 i = 1; % index for arrivals
0041 j = 1; % index for departures
0042 t = 0; % system time
0043
0044 while (t < A_max)
0045     t = min(A(i), D(j));
0046     if (t == A(i))
0047         N(k) = N(k-1) + 1;
0048         T(k) = t;
0049         k = k + 1;
0050         i = i + 1; % get next arrival
0051     else % departure occurs
0052         N(k) = N(k-1) - 1;
0053         T(k) = t;
0054         k = k + 1;
0055         j = j + 1; % get next departure
0056     end
0057 end
0058 %
0059 % record remaining departure instants
0060 for i=j:length(D)
0061     %
0062     N(k) = N(k-1) - 1;
0063     T(k) = t;
0064     k = k + 1;
0065 end
0066
0067 k = k - 1; % decrement k to get real size of N and T
0068 %
0069 % compute means
0070 MeanW = mean(W);
0071 T_Intervales = T(2:k)-T(1:k-1);
0072 MeanN = sum(N(1:k-1).*T_Intervales) / T(k);
0073 IdleDurationsIndex = find(N(1:k-1) == 0);
0074 Utilization = sum(T_Intervales(IdleDurationsIndex))/T(k);
0075 %

```

8/20/2005

Dr. Ashraf S.

Example: Queueing System – cont'd

```

0076 % Display results
0077 fprintf('Block #: '); fprintf('%3d ', [1:length(A)]); fprintf('\n');
0078 fprintf('Arrivals: '); fprintf('%3d ', A); fprintf('\n');
0079 fprintf('Errors: '); fprintf('%3d ', Errors); fprintf('\n');
0080 fprintf('Service: '); fprintf('%3d ', S); fprintf('\n');
0081 fprintf('Departs: '); fprintf('%3d ', D); fprintf('\n');
0082 fprintf('Waiting: '); fprintf('%3d ', W); fprintf('\n');
0083 fprintf('\n');
0084 fprintf('Average no of customers in system = %7.3f\n', MeanN);
0085 fprintf('Average customer waiting time = %7.3f microsec\n', MeanW);
0086 fprintf('Maximum simulation time = %7.3f microsec\n', T(k));
0087 fprintf('Duration server busy = %7.3f microsec\n', ...
sum(T_Intervales(IdleDurationsIndex)));
0088
0089 fprintf('Server utilisation = %7.4f\n', Utilization);
0090 fprintf('Server idle = %7.4f\n', 1.0-Utilization);
0091 %
0092 % Plot results
0093 figure(1)
0094 h = stairs(T, N); grid
0095 set(h, 'LineWidth', 3);
0096 xlabel('Time');
0097 ylabel('No of customers in system, N(t)');
0098
0099 figure(2);
0100 [AT, AA] = stairs(A, cumsum(ones(size(A))));
0101 [DT, DD] = stairs(D, cumsum(ones(size(D))));
0102 [NT, NN] = stairs(T, N);
0103 h = plot(AT, AA, '-', DT, DD, '--r', NT, NN, '-.-'); grid
0104 set(h, 'LineWidth', 3);
0105 title('Queue system simulation');
0106 ylabel('No of customers');
0107 xlabel('Time');
0108 legend('A(t)', 'D(t)', 'N(t)', 0);
0109
0110 figure(3);
0111 h = stem(W); grid
0112 set(h, 'LineWidth', 3);
0113 ylabel('Waiting time');
0114 xlabel('Customer index');
0115 LegendStr = ['MeanW = ' num2str(MeanW)];
0116 legend(LegendStr, 0);

```

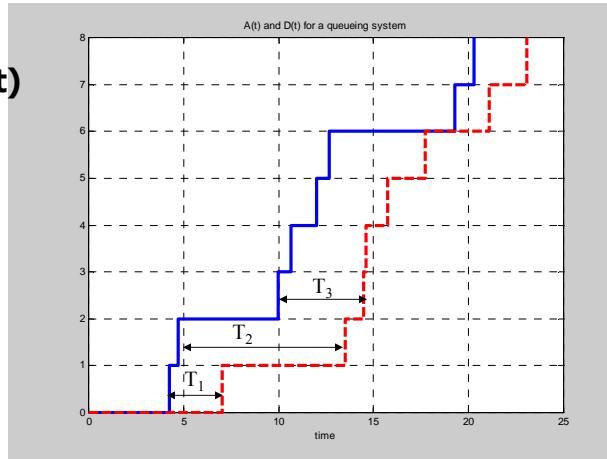
8/20/2005

Dr. Ashraf S. Hasan Mahmoud

8

Number of Customers in System

- **Blue curve: $A(t)$**
- **Red curve: $D(t)$**
- **Total time spent in the system for all customers = area in between two curves**



8/20/2005

Dr. Ashraf S. Hasan Mahmoud

9

Little's Formula

- **Little's formula:**
$$E[n] = \lambda E[r]$$

Holds for many service disciplines and for systems with arbitrary number of servers. It holds for many interpretations of the system as well

8/20/2005

Dr. Ashraf S. Hasan Mahmoud

10

Example 1:

- **Problem:** Let $n_s(t)$ be the number of customers being served at time t , and let τ denote the service time. If we designate the set of servers to be the "system" then Little's formula becomes:

$$E[n_s] = \lambda E[s]$$

Where $E[n_s]$ is the average number of busy servers for a system in the steady state.

Example 1: cont'd

Note: for a single server $n_s(t)$ can be either 0 or 1 $\rightarrow E[n_s]$ represents the portion of time the server is busy. If $p_0 = \text{Prob}[n_s(t) = 0]$, then we have

$$1 - p_0 = E[n_s] = \lambda E[s], \text{ Or} \\ p_0 = 1 - \lambda E[s]$$

The quantity $\lambda E[s]$ is defined as the utilization, U , for a single server. Usually, it is given the symbol ρ

$$\rho = \lambda E[s]$$

For a c -server system, we define the utilization (the fraction of busy servers) to be

$$\rho = \lambda E[s] / c$$

The M/M/1 Queue

- **Consider a single server system where customers arrive according to a Poisson process of rate λ**
 - **→ inter-arrival times are iid exponential r.v. with mean $E[\tau] = 1/\lambda$**
- **Assume the service times are iid exponential r.v. with mean $E[s] = 1/\mu$**
- **Assume the inter-arrival times and service times are independent**
- **Assume the system can accommodate unlimited number of customers**

8/20/2005

Dr. Ashraf S. Hasan Mahmoud

13

The M/M/1 Queue – cont'd

- **What is the steady state pmf of $n(t)$, the number of customers in the system?**
- **What is the PDF of r , the total customer delay in the system or the response time (as used in the textbook)?**

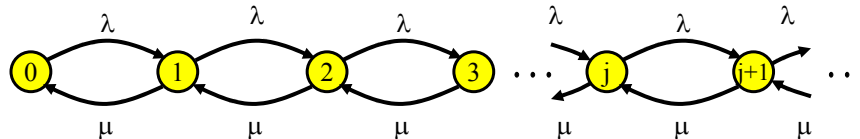
8/20/2005

Dr. Ashraf S. Hasan Mahmoud

14

The M/M/1 Queue – cont'd

- Consider the transition rate diagram for M/M/1 system



- **Note:**
 - System state – number of customers in systems
 - λ is rate of customer arrivals
 - μ is rate of customer departure

8/20/2005

Dr. Ashraf S. Hasan Mahmoud

15

The M/M/1 Queue – Distribution of Number of Customers

- Writing the global balance equations for this Markov chain and solving for $\text{Prob}[n(t) = j]$, yields (refer to previous example)

$$p_j = \text{Prob}[n(t) = j] \\ = (1 - \rho)\rho^j \quad \text{for } j=0,1,2, \dots, \\ \text{for } \rho = \lambda/\mu < 1$$

Note that for $\rho = 1 \rightarrow$ arrival rate $\lambda =$ service rate μ

8/20/2005

Dr. Ashraf S. Hasan Mahmoud

16

The M/M/1 Queue - Expected Number of Customers

- **The mean number of customers (in the whole system = buffer + server) is given by**

$$E[n] = \sum_j j \text{ Prob}[N(t) = j]$$

$$= \rho / (1 - \rho) \quad \text{for } \rho < 1$$

- **You can also show that the variance is equal to**

$$\text{Var}[n] = \sum_j (j - E[n])^2 \text{ Prob}[n(t) = j]$$

$$= \rho / (1 - \rho)^2 \quad \text{for } \rho < 1$$

8/20/2005

Dr. Ashraf S. Hasan Mahmoud

17

The M/M/1 Queue - Mean Customer Delay

- **The mean total response time for the system, $E[r]$, is found using Little's formula**

$$\begin{aligned} E[r] &= E[n] / \lambda \\ &= (1 / \mu) / (1 - \rho) \\ &= 1 / (\mu - \lambda) \end{aligned}$$

- **Actually, the response time CDF can be shown to be**

$$F(r) = 1 - e^{-r\mu(1-\rho)} \quad r \geq 0$$

- **Can you find $F^{-1}(r)$ for this CDF?**

8/20/2005

Dr. Ashraf S. Hasan Mahmoud

18

The M/M/1 Queue – Mean Queueing Time

- The mean waiting time in queue, $E[w]$, is given by

$$\begin{aligned} E[w] &= E[r] - E[s] \\ &= \rho / (1 - \rho) E[s] \end{aligned}$$

- The CDF for the waiting time can be shown to be

$$F(w) = 1 - \rho e^{-w\mu(1-\rho)} \quad w \geq 0$$

- Can you find $F^{-1}(w)$ for this CDF?

8/20/2005

Dr. Ashraf S. Hasan Mahmoud

19

The M/M/1 Queue – Mean Number in Queue

- Again we employ Little's formula:

$$\begin{aligned} E[n_q] &= \lambda E[w] \\ &= \rho^2 / (1 - \rho) \end{aligned}$$

Another way of finding $E[n_q]$ is using
 $E[n_q] = \sum_{n=1}^{\infty} (n-1)\rho^n = \sum_{n=1}^{\infty} (n-1)(1-\rho)\rho^n$

Remember:

$$\text{server utilization } \rho = \lambda / \mu = 1 - p_0$$

All previous quantities $E[n]$, $E[r]$, $E[w]$, and $E[n_q] \rightarrow \infty$ as $\rho \rightarrow 1$

8/20/2005

Dr. Ashraf S. Hasan Mahmoud

20

Busy Period

- **Def: The time interval between two successive idle intervals**

- **Mean busy period = $\frac{1}{\mu (1 - \rho)}$**

- **The textbook provides lots of other formulas in regard to the busy period.**

8/20/2005

Dr. Ashraf S. Hasan Mahmoud

21

Scaling Effect for M/M/1 Queues

- **Consider a queue of arrival rate λ whose service rate is μ**
 - $\rho = \lambda/\mu$,
 - **The expected delay $E[r]$ is given by**
$$E[r] = (1/\mu) / (1 - \rho)$$
- **If the arrival rate increases by a factor of K , then we either**
 - 1. Have K queueing systems, each with a server of rate μ**
 - 2. Have one queueing system with a server of rate $K \mu$**
- **Which of the two options will perform better?**

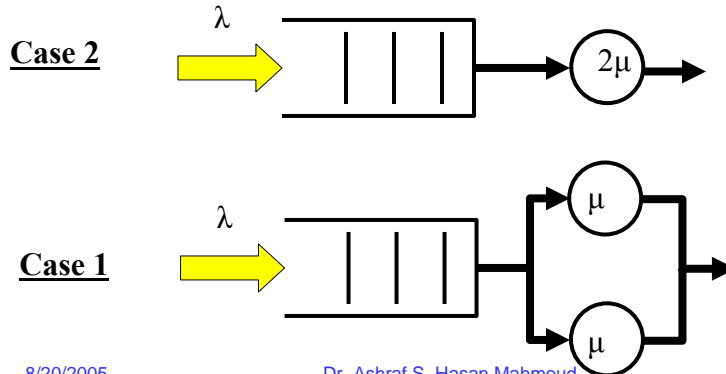
8/20/2005

Dr. Ashraf S. Hasan Mahmoud

22

Scaling Effect for M/M/1 Queues – cont'd

- **Example: $K = 2$: M/M/1 and M/M/2 systems with the same arrival rate and the same maximum processing rate**



8/20/2005

Dr. Ashraf S. Hasan Mahmoud

23

Scaling Effect for M/M/1 Queues – cont'd

- **Case 1: K queueing systems**
 - Identical systems
 - $E[r]$ is the same for all – $E[r] = (1/\mu) / (1-\rho)$
- **Case 2: 1 queueing system with server of rate $K\mu$**
 - ρ for this system = $(K\lambda) / (K\mu) = \lambda/\mu$ – same as the original system
 - $E[r'] = (1/(K\mu)) / (1-\rho) = (1/K) E[r]$
- **Therefore, the second option will provide a less total delay figure – significant delay performance improvement!**

8/20/2005

Dr. Ashraf S. Hasan Mahmoud

24

Example 31.1

- **Problem:** On a network gateway, measurements show that the packets arrive at a mean rate of 125 packets per second (pps) and the gateway takes about 2 milliseconds to forward them. Using an M/M/1 model:
 - Analyze the gateway.
 - What is the probability of buffer overflow had only 13 buffers?
 - How many buffers do we need to keep packet loss below one packet per million?

8/20/2005

Dr. Ashraf S. Hasan Mahmoud

25

Example 31.1 – cont'd

- **Solution:**

A) Analyzing the gateway:

Arrival rate, $\lambda = 125$ pps

Service rate, $\mu = 1/0.002 = 500$ pps

→ Gateway utilization, $\rho = \lambda / \mu = 0.25$

Prob[n packets in gateway] = $p_n = (1 - \rho)\rho^n$
 $= (0.75)(0.25)^n$

Mean # of packets in gateway, $E[n] = \rho / (1 - \rho)$
 $= 0.25 / 0.75 = 0.33$

Mean time in gateway, $E[r] = 1 / (\mu - \lambda) = 1 / (500 - 125) = 2.66$ msec

B) Prob[buffer overflow] = Prob[more than 13 packets in gateway]

$$\begin{aligned}
 &= p_{14} + p_{15} + p_{16} + \dots \\
 &= (1 - \rho)\{\rho^{14} + \rho^{15} + \rho^{16} + \dots\} \\
 &= (1 - \rho)\rho^{14}\{1 + \rho + \rho^2 + \dots\} \\
 &= \rho^{14} = 3.7 \times 10^{-9}
 \end{aligned}$$

C) To limit the probability of loss to less than 10^{-6} → Using (B)

Prob[buffer overflow] $\leq 10^{-6}$

→ Prob[more than n packets in gateway] $\leq 10^{-6}$

→ $\rho^{n+1} \leq 10^{-6}$ → $n \leq \log(10^{-6}) / \log(\rho) - 1 = 9.97 - 1$ → choose $n = 9$ buffers

Note –

(1) the solution in the textbook has typos in parts (B) and (C)
 (2) The solution for parts (B) and (C) are approximate – the more accurate model should be M/M/1/B – see the solution for this example after the M/M/1/B material

8/20/2005

Dr. Ashraf S. Hasan Mahmoud

26

M/M/1/B – Finite Capacity Queue

- Consider an M/M/1 with finite capacity $B < \infty$
- For this queue – there can be at most B customers in the system
 - 1 being served
 - $B-1$ waiting
- A customer arriving while the system has B customers is **BLOCKED** (does not wait)!

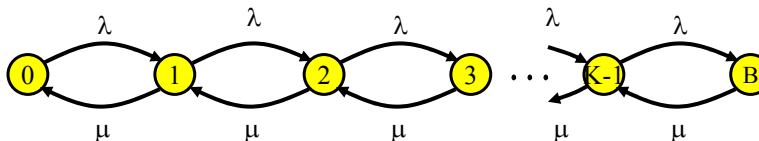
8/20/2005

Dr. Ashraf S. Hasan Mahmoud

27

M/M/1/B – Finite Capacity Queue – cont'd

- Transition rate diagram for this queueing system is given by:
 - $n(t)$ - A continuous-time Markov chain which takes on the values from the set $\{0, 1, \dots, B\}$



8/20/2005

Dr. Ashraf S. Hasan Mahmoud

28

M/M/1/B – Finite Capacity Queue – cont'd

- **The global balance equations:**

$$\begin{aligned} \lambda p_0 &= \mu p_1 \\ (\lambda + \mu)p_j &= \lambda p_{j-1} + \mu p_{j+1} \quad \text{for } j=1, 2, \dots, B-1 \\ \mu p_B &= \lambda p_{B-1} \end{aligned}$$

$$\begin{aligned} \rightarrow \text{Prob}[n(t) = j] &= p_j & j=0, 1, \dots, B; \rho < 1 \\ &= (1-\rho) \rho^j / (1-\rho^{B+1}) \end{aligned}$$

When $\rho = 1$, $p_j = 1/(B+1)$ (all states are equiprobable)

Will this system become unstable for $\rho = 1$? Why?

- **Two important numbers: p_0 and p_B**
 - p_0 is the probability of the server being idle – $p_0 = (1-\rho)/(1-\rho^{B+1})$
 - p_B is the probability of an arrival being blocked (or system overflow) – $p_B = (1-\rho) \rho^B / (1-\rho^{B+1})$

M/M/1/B – Mean Number of Customers

- **Mean number of customers, $E[n]$ is given by:**

$$\begin{aligned} E[n] &= \sum_{j=0}^B j \text{Pr}[n(t) = j] \\ &= \begin{cases} \frac{\rho}{1-\rho} - \frac{(B+1)\rho^{B+1}}{1-\rho^{B+1}} & \rho < 1 \\ B/2 & \rho = 1 \end{cases} \end{aligned}$$

M/M/1/B – Blocking Rate

- **A customer arriving while the system is in state K is BLOCKED (does not wait)!**
- **Therefore, rate of blocking, λ_b is given by**

$$\lambda_b = \lambda p_B$$

- **The actual arrival rate into the system is λ_a given**

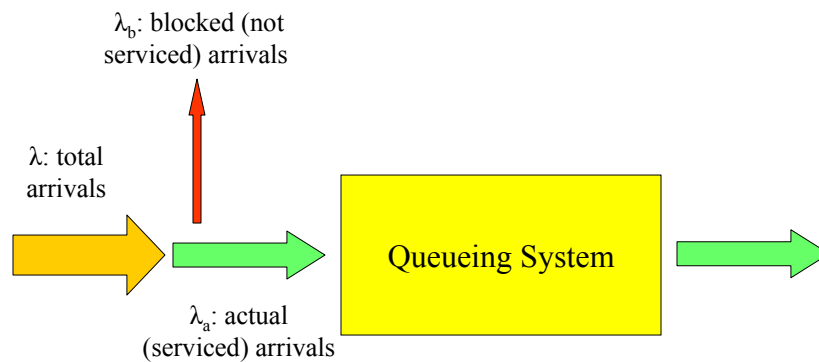
$$\begin{aligned}\lambda_a &= \lambda - \lambda_b \\ &= \lambda(1 - p_B)\end{aligned}$$

8/20/2005

Dr. Ashraf S. Hasan Mahmoud

31

M/M/1/B – Blocking Rate – cont'd



8/20/2005

Dr. Ashraf S. Hasan Mahmoud

32

M/M/1/B – Mean Response Time

- **The mean total response time, $E[r]$ is given by**

$$E[r] = E[n] / \lambda_a$$

8/20/2005

Dr. Ashraf S. Hasan Mahmoud

33

Example 31.1 – Redoing Part (B) and (C)

- **Problem: On a network gateway, measurements show that the packets arrive at a mean rate of 125 packets per second (pps) and the gateway takes about 2 milliseconds to forward them.**
 - Analyze the gateway.**
 - What is the probability of buffer overflow had only 13 buffers?**
 - How many buffers do we need to keep packet loss below one packet per million?**

8/20/2005

Dr. Ashraf S. Hasan Mahmoud

34

Example 31.1 – cont'd

- **Solution:**

B) Prob[buffer overflow] = p_{13}

$$= (1 - \rho) \rho^{13} / (1 - \rho^{13+1})$$

$$= 1.1 \times 10^{-8}$$

C) To limit the probability of loss to less than 10^{-6} → Using (B)

$$\text{Prob[buffer overflow]} \leq 10^{-6}$$

$$\rightarrow p_B \leq 10^{-6}$$

$$\rightarrow (1 - \rho) \rho^B / (1 - \rho^{B+1}) \leq 10^{-6}$$

$$\rightarrow \rho^B / (1 - \rho^{B+1}) \leq 10^{-6} / (1 - \rho) = 1.33 \times 10^{-6}$$

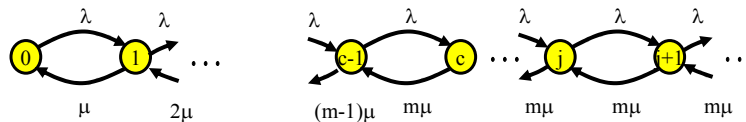
The above can be solved numerically for the value of B

$$\rightarrow B \geq 9.6 = 10 \text{ buffers}$$

35

Multi-Server Systems: M/M/m

- **The transition rate diagram for a multi-server M/M/m queue is as follows:**
 - **Departure rate = $k\mu$ when k servers are busy**



8/20/2005

Dr. Ashraf S. Hasan Mahmoud

36

Multi-Server Systems: M/M/m – cont'd

- **Writing the global balance equations:**

$$\lambda p_0 = \mu p_1$$

$$j\mu p_j = \lambda p_{j-1} \quad \text{for } j=1, 2, \dots, m$$

$$m\mu p_j = \lambda p_{j-1} \quad \text{for } j= m, m+1, \dots$$

→

$$p_j = a^j / j! p_0 \quad (\text{for } j=1, 2, \dots, m-1) \text{ and}$$

$$p_j = \rho^{j-m} / m! a^m p_0 \quad (\text{for } j=m, m+1, \dots)$$

where $a = \lambda/\mu$ and $\rho = a/m$

8/20/2005

Dr. Ashraf S. Hasan Mahmoud

37

Multi-Server Systems: M/M/m – cont'd

- **To find p_0 , we resort to the fact that $\sum p_j = 1$**

→

$$p_0 = \left\{ \sum_{j=0}^{m-1} \frac{a^j}{j!} + \frac{a^m}{m!} \frac{1}{1-\rho} \right\}^{-1}$$

Two important quantities:
- p_0 , and
- $\text{Prob}[W > 0]$

The probability that an arriving customer has to wait

$$\text{Prob}[W > 0] = \text{Prob}[N \geq m]$$

$$= p_m + p_{m+1} + p_{m+2} + \dots$$

$$= p_m / (1 - \rho)$$

$$= p_0 a^m / \{m! (1 - \rho)\}$$

In the textbook, this quantity is denoted by

Erlang-C formula

8/20/2005

Dr. Ashraf S. Hasan Mahmoud

38

Multi-Server Systems: M/M/m – cont'd

- **The mean number of customers in queue (waiting):**

$$\begin{aligned} E[n_q] &= \sum_{j=m}^{\infty} (j-m) \Pr[n(t) = j] \\ &= \sum_{j=m}^{\infty} (j-m) \rho^{j-m} p_m \\ &= \frac{\rho}{(1-\rho)^2} p_m \\ &= \frac{\rho}{1-\rho} \Pr[W > 0] \end{aligned}$$

8/20/2005

Dr. Ashraf S. Hasan Mahmoud

39

Multi-Server Systems: M/M/m – cont'd

- **The mean waiting time in queue:**

$$E[w] = E[n_q] / \lambda$$

- **The mean total response time in system:**

$$\begin{aligned} E[r] &= E[w] + E[s] \\ &= E[w] + 1 / \mu \end{aligned}$$

- **The mean number of customers in system:**

$$E[n] = \lambda E[r]$$

$$= E[n_q] + a$$

Why?

- **In addition to the formulas above the textbook gives the corresponding formulas for the variances too.**

8/20/2005

Dr. Ashraf S. Hasan Mahmoud

40

Example 2:

- **A company has a system with four private telephone lines connecting two of its sites. Suppose that requests for these lines arrive according to a Poisson process at rate of one call every 2 minutes, and suppose that call durations are exponentially distributed with mean 4 minutes. When all lines are busy, the system delays (i.e. queues) call requests until a line becomes available. Find the probability of having to wait for a line.**

8/20/2005

Dr. Ashraf S. Hasan Mahmoud

41

Example 2: cont'd

- **Solution:**

$$\lambda = 1/2, 1/\mu = 4, m = 4 \rightarrow a = \lambda / \mu = 2$$
$$\rightarrow \rho = a/m = 1/2$$

$$p_0 = \{1 + 2 + 2^2/2! + 2^3/3! + 2^4/4! (1/(1-\rho))\}^{-1}$$
$$= 3/23$$

$$p_m = a^m/m! p_0$$
$$= 2^4/4! \times 3/23$$

$$\text{Prob}[W > 0] = p_m/(1-\rho)$$
$$= 2^4/4! \times 3/23 \times 1/(1-1/2)$$
$$= 4/23$$
$$\approx 0.17$$

8/20/2005

Dr. Ashraf S. Hasan Mahmoud

42

Waiting/Response Time Distribution for M/M/m

- **See textbook**

$$F(w) = 1 - \frac{P_c}{1 - \rho} e^{-m\mu(1-\rho)w} \quad w > 0$$

- **The q-percentile can be computed as follows:**

$$\begin{aligned} w_q &= \max \left\{ 0, \frac{1}{m\mu(1-\rho)} \ln \left(\frac{100 * \Pr[W > 0]}{100 - q} \right) \right\} \\ &= \max \left\{ 0, \frac{E[w]}{\Pr[W > 0]} \ln \left(\frac{100 * \Pr[W > 0]}{100 - q} \right) \right\} \end{aligned}$$

8/20/2005

Dr. Ashraf S. Hasan Mahmoud

43

Example 13.2:

- **Problem:** Students arrive at the university computer center in a Poisson manner at an average rate of 10 per hour. Each student spends an average of 20 minutes at the terminal, and the time can be assumed to be exponentially distributed. The center currently has five terminals. Some students have been complaining that waiting times are too long.
- **Analyze the center using a queueing model.**

8/20/2005

Dr. Ashraf S. Hasan Mahmoud

44

$$p_0 = \left\{ \sum_{j=0}^{m-1} \frac{a^j}{j!} + \frac{a^m}{m!} \frac{1}{1-\rho} \right\}^{-1}$$

Example 13.2: cont'd

• **Solution:**

The center model: M/M/5

Arrival rate, $\lambda = 1/6$ student/min

Service rate, $\mu = 1/20$ student/min

→ Center utilization: $a = \lambda / \mu = 3.3333$

$\rho = a/m = 3.3333/5 = 0.6667 \leftarrow$ avg server

utilization

$$p_0 = \{1 + 10/3 + (10/3)^2/2! + (10/3)^3/3! + (10/3)^4/4! + (10/3)^5/5!(1-2/3)\}^{-1}$$

$$= \{31.4938\}^{-1} = 0.0318$$

$$p_m = a^m/m! p_0$$

$$= (10/3)^5 / 5! * 0.0318$$

$$= 0.1091$$

$$\text{Prob}[W>0] = p_m/(1-\rho)$$

$$= 0.1091/(1-2/3)$$

$$= 0.3271 \leftarrow \text{This is the probability that a student has to wait!}$$

8/20/2005

Dr. Ashraf S. Hasan Mahmoud

45

$$E[n_q] = \frac{\rho}{1-\rho} \text{Pr}[W > 0]$$

Example 13.2: cont'd

• **Solution:**

Avg # of students waiting in the center, $E[n_q]$

$$E[n_q] = \rho/(1-\rho) \text{Prob}[W>0]$$

$$= (2/3)/(1-2/3) * 0.3271$$

$$= 0.6543 \text{ students}$$

Avg waiting time for a student, $E[w]$

$$E[w] = E[n_q] / \lambda$$

$$= 0.6543 / (1/6)$$

$$= 3.9259 \approx 4 \text{ minutes}$$

Avg time spent in the centre, $E[r]$

$$E[r] = E[w] + E[s]$$

$$= 4 + 1/(1/20)$$

$$= 24 \text{ minutes}$$

Avg/variance # of students in the center, $E[n]$

$$E[n] = \lambda E[r]$$

$$= (1/6) * 24 = 4 \text{ students}$$

$\text{Var}[n] = \dots = 479 \text{ students}^2$

8/20/2005

Dr. Ashraf S. Hasan Mahmoud

46

$w_q = \max\left\{0, \frac{1}{m\mu(1-\rho)} \ln\left(\frac{100 \cdot \Pr[W > 0]}{100 - q}\right)\right\}$

Example 13.2: cont'd

- **Solution:**

Avg # of students using the terminals, $E[n_s]$

$$\begin{aligned} E[n_s] &= E[n] - E[nq] \\ &= 4 - 0.6543 \\ &= 3.35 \text{ students} \end{aligned}$$

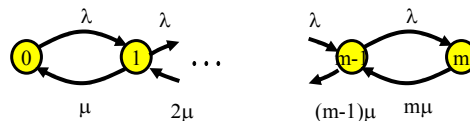
The 90-percentile of the waiting time, $w_{0.9}$

$$\begin{aligned} w_{0.9} &= \max\{0, 4/0.3271 \ln(10 \cdot 0.3271)\} \\ &= 14 \text{ minutes} \end{aligned}$$

THUS, only 10% of the students have to wait more than 14 minutes!!

Multi-Server Systems: M/M/m/m

- **The transition rate diagram for a multi-server with no waiting room (M/M/m/m) queue is as follows:**
 - **Departure rate = $k\mu$ when k servers are busy**



PMF for Number of Customers for M/M/m/m

- **Writing the global balance equations, one can show:**

$$p_j = a^j / j! p_0 \quad (\text{for } j=0, 1, \dots, m)$$

where $a = \lambda / \mu$ (the offered load)

- **To find p_0 , we resort to the fact that $\sum p_j = 1$**

$$p_0 = \left\{ \sum_{j=0}^m \frac{a^j}{j!} \right\}^{-1}$$

Erlang-B Formula

- **Erlang-B formula is defined as the probability that all servers are busy:**

$$\begin{aligned} \Pr[n = m] &= p_m \\ &= \frac{a^m / m!}{1 + a + a^2 / 2! + \dots + a^m / m!} \end{aligned}$$

Expected Number of customers in M/M/m/m

- **The actual arrival rate *into* the system:**

$$\lambda_a = \lambda(1 - p_m)$$

- **Average total delay figure:**

$$E[r] = E[s] = 1 / \mu$$

Why?

- **Average number of customers:**

$$E[n] = \lambda_a E[s] = \lambda_a / \mu$$