

King Fahd University of Petroleum & Minerals Computer Engineering Dept

COE 540 – Computer Networks
Term 081
Dr. Ashraf S. Hasan Mahmoud
Rm 22-148-3
Ext. 1724
Email: ashraf@kfupm.edu.sa

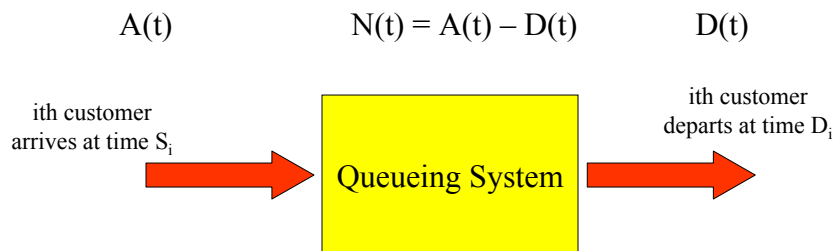
1/6/2009

Dr. Ashraf S. Hasan Mahmoud

1

Queuing Model

- Consider the following system:



$$T_i = D_i - A_i$$

$$W_i = T_i - S_i \\ = D_i - A_i - S_i$$

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

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

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

T_i – duration of time spent in system for i th customer

W_i – duration of time spent waiting for service for i th customer

2

Example 1: 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?

1/6/2009

Dr. Ashraf S. Hasan Mahmoud

3

Example 1: 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 |

1/6/2009

Dr. Ashraf S. Hasan Mahmoud

4

Example 1: Queueing System – cont'd

Solution:

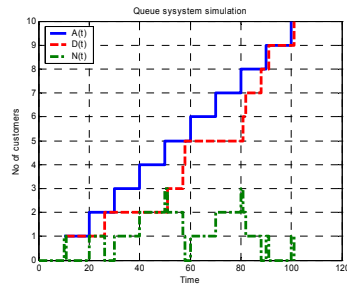
Using the previous results and knowing that

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

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)



1/6/2009

Dr. Ashraf S. Hasan M

Example 1: 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     t = D(i);
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_Intervals = T(2:k)-T(1:k-1);
0072 MeanN = sum(N(1:k-1))*T_Intervals / T(k);
0073 IdleDurationsIndex = find(N(1:k-1) == 0);
0074 Utilization = sum(T_Intervals(IdleDurationsIndex))/T(k);
0075 %
    
```

1/6/2009

Dr. Ashraf S.

Example 1: 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\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', ...
0088 sum(T_Intervales(IdleDurationsIndex)));
0089 fprintf('Server utilization = %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, '-l', 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);

```

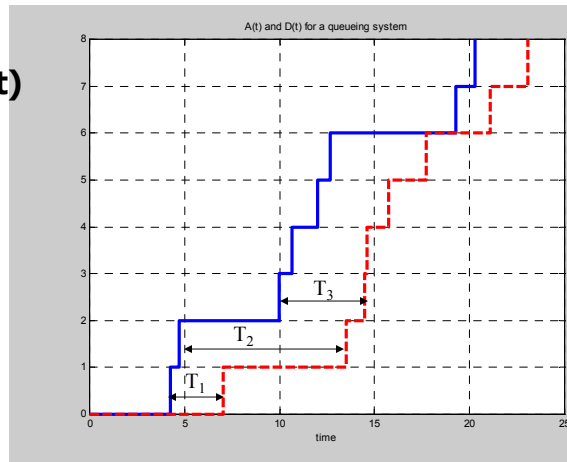
1/6/2009

Dr. Ashraf S. Hasan Mahmoud

7

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**



1/6/2009

Dr. Ashraf S. Hasan Mahmoud

8

Little's Formula

- Little's formula:

$$E[N] = \lambda E[T]$$

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

Example 2:

- **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" m then Little's formula becomes:

$$E[N_s] = \lambda E[\tau]$$

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

Example 2: 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

$$1 - p_0 = E[N_s] = \lambda E[\tau], \text{ Or} \\ p_0 = 1 - \lambda E[\tau]$$

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

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

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

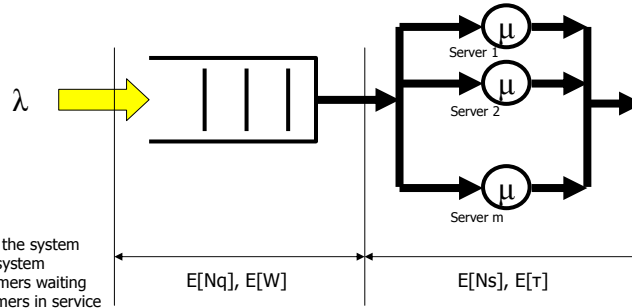
$$\rho = \lambda E[\tau] / c$$

Example 3: Applications on Little's Formula

**Refer to the slides for Dr. Waheed.
The slides have 8 good examples!**

Queue System and Parameters

- Queueing system with m servers
 - When $m = 1$ – single server system
- Input: arrival statistics (rate λ), service statistics (rate μ), number of customers (m), buffer size
- Output: $E[N]$, $E[T]$, $E[Nq]$, $E[W]$, $\text{Prob}[\text{buffer size} = x]$, $\text{Prob}[W < w]$, etc.



$E[N]$ = mean # of customers in the system
 $E[T]$ = mean time spent in the system
 $E[Nq]$ = mean number of customers waiting
 $E[Ns]$ = mean number of customers in service
 $E[W]$ = mean waiting time for a customer
 $E[\tau]$ = mean service time for a customer

$$\begin{aligned}
 E[N], E[T] \\
 E[N] &= E[Nq] + E[Ns], \\
 E[T] &= E[W] + E[\tau]
 \end{aligned}$$

1/6/2009

Dr. Ashraf S. Hasan Mahmoud

13

The M/M/1 Queue

- Consider m -server system where customers arrive according to a Poisson process of rate λ
 - \rightarrow inter-arrival times are iid exponential r.v. with mean $1/\lambda$
- Assume the service times are iid exponential r.v. with mean $1/\mu$
- Assume the inter-arrival times and service times are independent
- Assume the system can accommodate unlimited number of customers

1/6/2009

Dr. Ashraf S. Hasan Mahmoud

14

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 T , the total customer delay in the system?

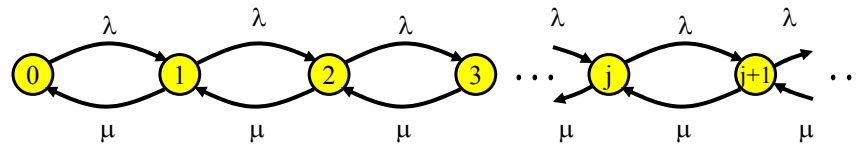
1/6/2009

Dr. Ashraf S. Hasan Mahmoud

15

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

1/6/2009

Dr. Ashraf S. Hasan Mahmoud

16

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)

$$p_j = \text{Prob}[N(t) = j] \\ = (1-\rho)\rho^j$$

for $\rho = \lambda/\mu < 1$

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

1/6/2009

Dr. Ashraf S. Hasan Mahmoud

17

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

- The mean number of customer is given by

$$E[N] = \sum_j j \text{Prob}[N(t) = j] \\ = \rho / (1-\rho)$$

1/6/2009

Dr. Ashraf S. Hasan Mahmoud

18

The M/M/1 Queue – Mean Customer Delay

- The mean total customer delay in the system is found using Little's formula

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

1/6/2009

Dr. Ashraf S. Hasan Mahmoud

19

The M/M/1 Queue – Mean Queueing Time

- The mean waiting time in queue is given by

$$\begin{aligned}E[W] &= E[T] - E[\tau] \\ &= \rho / (1 - \rho) E[\tau]\end{aligned}$$

1/6/2009

Dr. Ashraf S. Hasan Mahmoud

20

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

- Again we employ Little's formula:

$$E[N_q] = \lambda E[W]$$

$$= \rho^2 / (1-\rho)$$

Remember:

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

All previous quantities $E[N]$, $E[T]$, $E[W]$, and $E[N_q] \rightarrow \infty$ as $\rho \rightarrow 1$

1/6/2009

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 $E[T]$ is given by
$$E[T] = (1/\mu) / (1-\rho)$$
- 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?

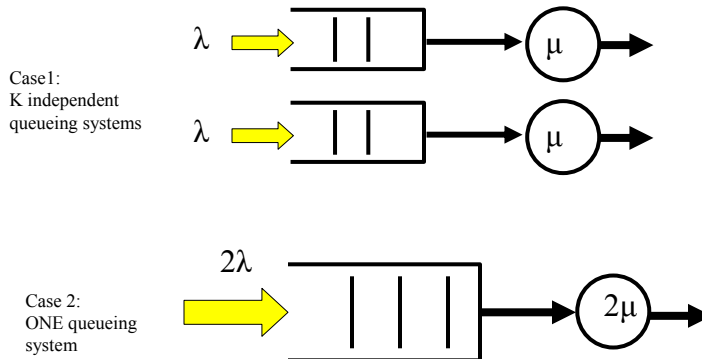
1/6/2009

Dr. Ashraf S. Hasan Mahmoud

22

Example 4: Scaling Effect for M/M/1 Queues

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



1/6/2009

Dr. Ashraf S. Hasan Mahmoud

23

Example 4: Scaling Effect for M/M/1 Queues – cont'd

- **Case 1: K queueing systems**
 - Identical systems
 - $E[T]$ is the same for all – $E[T] = (1/\mu) / (1-\rho)$
- **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
 - $E[T'] = (1/(K\mu)) / (1-\rho) = (1/K) E[T]$
- **Therefore, the second option will provide a less total delay figure – significant delay performance improvement!**

1/6/2009

Dr. Ashraf S. Hasan Mahmoud

24

M/M/1/K – Finite Capacity Queue

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

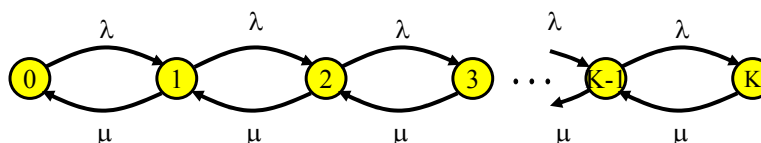
1/6/2009

Dr. Ashraf S. Hasan Mahmoud

25

M/M/1/K – 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 $\{0, 1, \dots, K\}$



1/6/2009

Dr. Ashraf S. Hasan Mahmoud

26

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

- The global balance equations:

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

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

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

1/6/2009

Dr. Ashraf S. Hasan Mahmoud

27

M/M/1/K – Mean Number of Customers

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

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

1/6/2009

Dr. Ashraf S. Hasan Mahmoud

28

M/M/1/K – 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**

$$\lambda_b = \lambda p_K$$

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

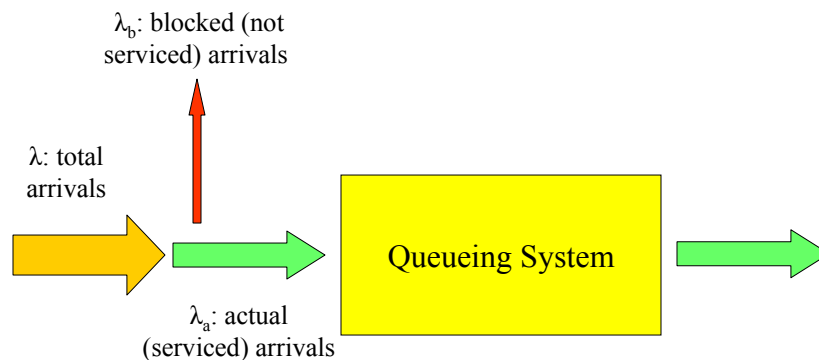
$$\begin{aligned}\lambda_a &= \lambda - \lambda_b \\ &= \lambda(1 - p_K)\end{aligned}$$

1/6/2009

Dr. Ashraf S. Hasan Mahmoud

29

M/M/1/K – Blocking Rate – cont'd



1/6/2009

Dr. Ashraf S. Hasan Mahmoud

30

M/M/1/K – Mean Delay

- The mean total delay $E[T]$ is given by

$$E[T] = E[N] / \lambda_a$$

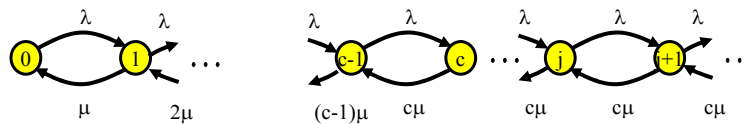
1/6/2009

Dr. Ashraf S. Hasan Mahmoud

31

Multi-Server Systems: M/M/c

- The transition rate diagram for a multi-server M/M/c queue is as follows:
 - Departure rate = $k\mu$ when k servers are busy
 - We can show that the service time for a customer finding k servers busy is exponentially distributed with mean $1/(k\mu)$



1/6/2009

Dr. Ashraf S. Hasan Mahmoud

32

Multi-Server Systems: M/M/c – cont'd

- Writing the global balance equations:

$$\begin{array}{ll} \lambda & p_0 = \mu p_1 \\ j\mu & p_j = \lambda p_{j-1} \quad \text{for } j=1, 2, \dots, c \\ c\mu & p_j = \lambda p_{j-1} \quad \text{for } j= c, c+1, \dots \end{array}$$

Note this distribution is the same as that for M/M/1 when you set c to 1.

→

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

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

- From this we note that the probability of system being in state c, p_c , is given by

$$p_c = a^c / c! p_0$$

1/6/2009

Dr. Ashraf S. Hasan Mahmoud

33

Multi-Server Systems: M/M/c – cont'd

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

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

- The probability that an arriving customer has to wait

$$\begin{aligned} \text{Prob}[W > 0] &= \text{Prob}[N \geq c] \\ &= p_c + p_{c+1} + p_{c+2} + \dots \\ &= p_c / (1-\rho) \end{aligned}$$

Erlang-C formula

Question: What is $\text{Prob}[W>0]$ for M/M/1 system?

1/6/2009

Dr. Ashraf S. Hasan Mahmoud

34

Multi-Server Systems: M/M/c – cont'd

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

$$\begin{aligned} E[N_q] &= \sum_{j=c}^{\infty} (j-c) \Pr[N(t) = j] \\ &= \sum_{j=c}^{\infty} (j-c) \rho^{j-c} p_c \\ &= \frac{\rho}{(1-\rho)^2} p_c \\ &= \frac{\rho}{1-\rho} \Pr[W > 0] \end{aligned}$$

1/6/2009

Dr. Ashraf S. Hasan Mahmoud

35

Multi-Server Systems: M/M/c – cont'd

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

$$E[W] = E[N_q] / \lambda$$

- **The mean total delay in system:**

$$\begin{aligned} E[T] &= E[W] + E[\tau] \\ &= E[W] + 1/\mu \end{aligned}$$

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

$$\begin{aligned} E[N] &= \lambda E[T] \\ &= E[N_q] + a \end{aligned}$$

Why?

1/6/2009

Dr. Ashraf S. Hasan Mahmoud

36

Example 5:

- 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.
- What is the average waiting time for an incoming call?

1/6/2009

Dr. Ashraf S. Hasan Mahmoud

37

Example 5: cont'd

- **Solution:**
 $\lambda = 1/2, 1/\mu = 4, c = 4 \rightarrow a = \lambda/\mu = 2$
 $\rightarrow \rho = a/c = 1/2$
 $p_0 = \{1 + 2 + 2^2/2! + 2^3/3! + 2^4/4! (1/(1-\rho))\}^{-1}$
 $= 3/23$
 $p_c = a^c/c! p_0$
 $= 2^4/4! \times 3/23$

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

(2) To find $E[W]$, find $E[N_q]$...
 $E[N_q] = \rho/(1-\rho) * \text{Prob}[W > 0] = 0.1739$
 $E[W] = E[N_q]/\lambda = 0.35 \text{ min}$

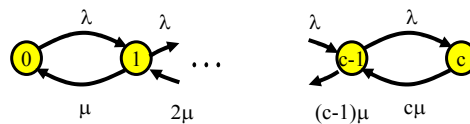
1/6/2009

Dr. Ashraf S. Hasan Mahmoud

38

Multi-Server Systems: M/M/c/c

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



1/6/2009

Dr. Ashraf S. Hasan Mahmoud

39

PMF for Number of Customers for M/M/c/c

- Writing the global balance equations, one can show:

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

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

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

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

1/6/2009

Dr. Ashraf S. Hasan Mahmoud

40

Erlang-B Formula

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

$$\Pr[N = c] = p_c = \frac{a^c / c!}{1 + a + a^2 / 2! + \dots + a^c / c!}$$

Expected Number of customers in M/M/c/c

- The actual arrival rate *into* the system:

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

- Average total delay figure:

$$E[T] = E[\tau]$$

Why?

- Average number of customers:

$$E[N] = \lambda_a E[\tau]$$

Example 6:

- 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 BLOCKS the incoming call and generates a busy signal.
- Find the probability of being blocked.

1/6/2009

Dr. Ashraf S. Hasan Mahmoud

43

Example 6:

- **Solution:**
 $\lambda = 1/2, 1/\mu = 4, c = 4 \rightarrow a = \lambda/\mu = 2$
 $\rightarrow \rho = a/c = 1/2$

$$p_c = \frac{a^c/c!}{1 + a + a^2/2! + a^3/3! + a^4/4!}$$
$$= \frac{2^4/4!}{1 + 2 + 2^2/2! + 2^3/3! + 2^4/4!} = 9.5\%$$

Therefore, the probability of being blocked is 0.095.

1/6/2009

Dr. Ashraf S. Hasan Mahmoud

44

M/G/1 Queues

- **Poisson arrival process (i.e. exponential r.v. interarrival times)**
- **Service time: general distribution $f_{\tau}(x)$**
 - For M/M/1, $f_{\tau}(x) = \mu e^{-\mu x}$ for $x > 0$
- **The state of the M/G/1 system at time t is specified by**
 1. **$N(t)$**
 2. **The remaining (residual) service time of the customer being served**

1/6/2009

Dr. Ashraf S. Hasan Mahmoud

45

The Residual Service Time

- **Mean residual time (see example and derivation in handout) is given by**

$$E[R] = \frac{E[\tau^2]}{2E[\tau]}$$

1/6/2009

Dr. Ashraf S. Hasan Mahmoud

46

Mean Waiting Time in M/G/1

- The waiting time of a customer is the sum of the residual service time R' of the customer (if any) found in service and the $N_q(t) = k-1$ service time of the customers (if any) found in queue

$$\begin{aligned}E[W] &= E[R'] + E[N_q] E[\tau] \\ &= E[R'] + \lambda E[W] E[\tau] \\ &= E[R'] + \rho E[W]\end{aligned}$$

1/6/2009

Dr. Ashraf S. Hasan Mahmoud

47

Mean Waiting Time in M/G/1 – cont'd

- But residual service time R' (as observed by an arriving customer) is either
 - 0 if the server is free
 - R if the server is busy
- Therefore, mean of R' is given by

$$\begin{aligned}E[R'] &= 0 \times \text{Pro}[N(t)=0] + E[R](1-\text{Pro}[N(t)=0]) \\ &= E[\tau^2]/(2E[\tau]) \times \rho \\ &= \lambda E[\tau^2]/2\end{aligned}$$

1/6/2009

Dr. Ashraf S. Hasan Mahmoud

48

Mean Waiting Time in M/G/1 – cont'd

- Substituting back, yields

$$E[W] = \frac{\lambda E[\tau^2]}{2(1-\rho)}$$

$$= \frac{\lambda(\delta_\tau^2 + E[\tau]^2)}{2(1-\rho)}$$

$$= \frac{\rho(1 + C_\tau^2)}{2(1-\rho)} E[\tau]$$

Remember:

$$- E[\tau^2] = \delta_\tau^2 + E[\tau]^2$$

$$- C_\tau^2 = \delta_\tau^2 / E[\tau]^2$$

Pollaczek-Khinchin (P-K)
Mean Value Formula

Mean Delay in M/G/1 – cont'd

- The mean waiting time, $E[T]$ is found by adding mean service time to $E[W]$:

$$E[T] = E[\tau] + E[W]$$

$$= E[\tau] + \frac{\rho(1 + C_\tau^2)}{2(1-\rho)} E[\tau]$$

Example 7:

- **Problem:** Compare $E[W]$ for M/M/1 and M/D/1 systems.

- **Answer:**

M/M/1: service time, τ , is exponential r.v. with parameter μ

$$\rightarrow E[\tau] = 1/\mu, E[\tau^2] = 2/\mu^2, \delta^2_{\tau} = 1/\mu^2, C^2_{\tau} = 1$$

M/D/1: service time, τ , is constant with value $\tau = 1/\mu$

$$\rightarrow E[t] = 1/\mu, E[\tau^2] = 1/\mu^2, \delta^2_{\tau} = 0, C^2_{\tau} = 0$$

1/6/2009

Dr. Ashraf S. Hasan Mahmoud

51

Example 7: cont'd

- **Answer:** cont'd

Substitute in P-K mean value formula

M/M/1:

$$E[W_{M/M/1}] = \frac{\lambda E[\tau^2]}{2(1-\rho)} = \frac{\rho}{(1-\rho)} E[\tau]$$

M/D/1:

$$E[W_{M/D/1}] = \frac{\lambda E[\tau^2]}{2(1-\rho)} = \frac{\rho}{2(1-\rho)} E[\tau]$$

$$= \frac{1}{2} E[W_{M/M/1}]$$

The waiting time in an M/D/1 queue is half of that of an M/M/1 system

1/6/2009

Dr. Ashraf S. Hasan Mahmoud

52

Example 8:

- **Problem:** Assume traffic is arriving at the input port of a router according to a Poisson arrival process of rate $\lambda = 100$ packets/sec. If the traffic distribution is as follows:
 - 30% of packets are 512 Bytes long,
 - 50% of packets are 1024 Bytes long,
 - 20% of packets are 4096 Bytes longIf the transmit speed of the router output port is 1.5 Mb/s
 - a) What is the average packet transmit time?
 - b) What is the average packet waiting time before transmit?
 - c) What is the average buffer size in the router?

1/6/2009

Dr. Ashraf S. Hasan Mahmoud

53

Example 8: cont'd

- **Solution:**
 - a) Average packet size,
 $E[L] = 0.3 \times 512 + 0.5 \times 1024 + 0.2 \times 4096$
 $= 1484.8$ Bytes
average transmit time = $E[L]/R = 1484.8 \times 8 / 1.5 \times 10^6 = 0.0079$ sec
 - b) $E[L^2] = 0.3 \times (512 \times 8)^2 + 0.5 \times (1024 \times 8)^2 + 0.2 \times (4096 \times 8)^2 = 2.5334 \times 10^8$ Bits²
 $E[\tau^2] = E[L^2]/R^2 = 1.1259 \times 10^{-4}$ sec²
 $\rho = \lambda E[\tau] = 0.7919$
 $E[W] = 0.5 \lambda E[\tau^2] / (1 - \rho)$
 $= 0.0271$ sec
 - c) $E[Nq] = \lambda E[W]$
 $= 2.705$ packet

1/6/2009

Dr. Ashraf S. Hasan Mahmoud

54