

# The G/M/1, G/G/1, G/G/m and M/G/m/m Queues

Copyright 2002, Sanjay K. Bose

1

## The G/M/1 Queue

- The G/M/1 queue is the dual of the M/G/1 queue where the arrival process is a general one but the service times are exponentially distributed.
- Service time distribution is exponential with parameter  $1/m$
- General Arrival Process with mean arrival rate  $I$ .  
Inter-arrival time is random with pdf  $a(t)$ , cdf  $A(t)$   
and L.T. of the pdf as  $L_A(s)$
- Total Traffic  $r = I/m$

Stability consideration require that  $r < 1$  for the queue to be at equilibrium

Copyright 2002, Sanjay K. Bose

2

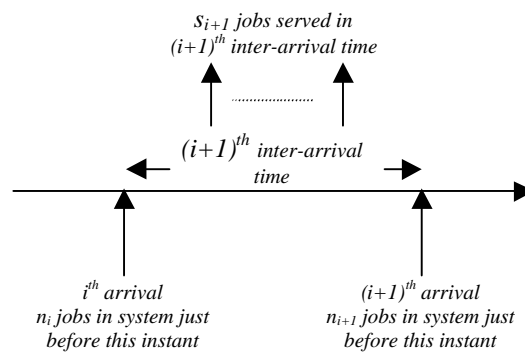
- For analyzing the G/M/1 queue using the Imbedded Markov Chain approach, the *imbedded points are chosen to be the arrival instants of jobs to the system*

- System State = Number in the system *immediately before an arrival instant*

$n_i$  = Number in the system just before the  $i^{\text{th}}$  arrival

$s_{i+1}$  = Number of jobs served between the  $i^{\text{th}}$  and the  $(i+1)^{\text{th}}$  arrivals

The sequence  $\{n_i\}$ ,  $i=1,2,\dots$  at the imbedded Markov points (i.e. just before arrival instants) forms a Markov Chain.



$$n_{i+1} = n_i + 1 - s_{i+1} \quad n_i = 0, 1, \dots, \infty \quad s_{i+1} \in n_i + 1 \quad (1)$$

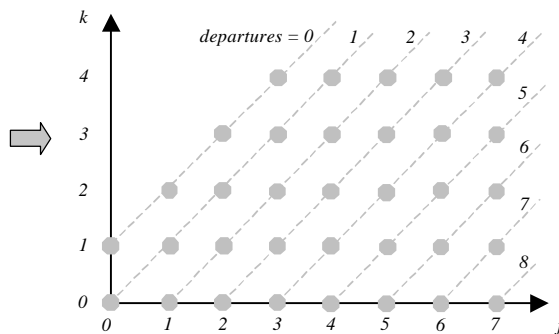
Consider the chain  $n_{i+1} = n_i + 1 - s_{i+1}$   $n_i = 0, 1, \dots, \infty$   $s_{i+1} \in n_i + 1$   
 at equilibrium  
 (i.e. for  $i \in \mathbb{N}$ )

One-Step Transition Probabilities  $\begin{cases} p_{jk} = P\{n_{i+1} = k | n_i = j\} \\ p_{jk} = 0 \text{ for } k > j + 1 \end{cases}$  (2)

Balance Equations  $p_k = \sum_{j=0}^{\infty} p_j p_{jk}$   $k = 0, 1, \dots, \infty$  (3)

Normalization Condition  $\sum_{k=0}^{\infty} p_k = 1$

Allowable state transitions ( $j \rightarrow k$ ) in the G/M/1 queue between successive imbedded points



The points shown are the ones for which the ( $j \rightarrow k$ ) transitions can occur. For all other points, the corresponding transitions cannot occur.

The diagonal lines correspond to a constant number of departures (0, 1, 2, ...) between the corresponding ( $j \rightarrow k$ ) transitions.

Let  $\mathbf{a}_n = P\{n \text{ departures in an inter-arrival time interval } | \text{ server is always busy during this inter-arrival time}\}$

$$\mathbf{a}_j = \int_{x=0}^{\infty} \frac{(\mathbf{m}x)^j}{j!} e^{-\mathbf{m}x} a(x) dx \quad j = 0, 1, \dots, \infty \quad (5)$$

**Justification:** *Since the service times are exponentially distributed, the number of departures in an inter-arrival time instant where the server is always busy will have the Poisson distribution.*

and 
$$\sum_{j=0}^{\infty} \mathbf{a}_j z^j = \int_0^{\infty} e^{-\mathbf{m}(1-z)} a(x) dx = L_A(\mathbf{m} - \mathbf{m}z)$$



The  $\mathbf{a}_j$ 's may be found as the coefficient of  $z^j$  in the series expansion of  $L_A(\mathbf{m} - \mathbf{m}z)$

Copyright 2002, Sanjay K. Bose

7

The Balance Equations of (3) then become

$$p_j = \mathbf{a}_0 p_{j-1} + \sum_{k=0}^{\infty} \mathbf{a}_{k+1} p_{j+k} \quad j = 1, \dots, \infty \quad (4)$$

Solutions to these, satisfying the normalization condition are -

$$p_j = (1 - \mathbf{s}) \mathbf{s}^j \quad j = 0, 1, \dots, \infty \quad (6)$$

where  $\mathbf{s}$  is a unique root of  $\mathbf{s} = L_A(\mathbf{m} - \mathbf{m}\mathbf{s})$

*It may be shown that if  $\mathbf{r} < 1$ , then there will always be a unique real solution for  $\mathbf{s} = L_A(\mathbf{m} - \mathbf{m}\mathbf{s})$  which will be in the range  $0 < \mathbf{s} < 1$ .*

Copyright 2002, Sanjay K. Bose

8

The solution may be verified by direct substitution, as given below -

$$(1-s)s^j = a_0(1-s)s^{j-1} + \sum_{k=0}^{\infty} a_{k+1}(1-s)s^{j+k}$$

$$s^j = a_0s^{j-1} + \sum_{k=1}^{\infty} a_k s^{j+k-1}$$

$$s = a_0 + \sum_{k=1}^{\infty} a_k s^k$$

$$= \sum_{k=0}^{\infty} a_k s^k = L_A(m - ms)$$

Copyright 2002, Sanjay K. Bose

9

- The state distribution  $\{p_j\}$   $j=0,1,\dots$  is found under equilibrium conditions at the time instants just before job arrivals to the system.

- It is also valid for the departure instants (just after a job leaves the system) as Kleinrock's principle is applicable to this system (i.e. the state changes are at most  $+1$  or  $-1$ ).

- It is **not valid** for arbitrary time instants (or ergodic, time-average results) since PASTA will not be applicable to the system (i.e. the arrival process is not Poisson).

Copyright 2002, Sanjay K. Bose

10

For a FCFS M/G/1 queue at equilibrium, the following queueing delay results may be obtained.

These may be derived using the fact that the service time distribution is exponential and hence memory less

$$W_q = \sum_{n=1}^{\infty} \frac{n}{m} (1-s) s^n = \frac{s}{m(1-s)}$$

$$f_{W_q}(t) = (1-s)d(t) + ms(1-s)e^{-m(1-s)t}$$

for  $t \geq 0$

Other parameters such as  $W$  and  $N_q$  and the distribution  $f_w(t)$  may also be obtained

The multi-server G/M/m queue may also be analyzed using a similar approach

### The G/G/1 Queue

We cannot analyse this queue exactly but there are useful bounds that have been developed for the waiting time in queue  $W_q$ . This can then be used to find bounds on  $W$ ,  $N$  and  $N_q$  in the usual fashion, i.e. Little's Result and  $W = W_q + \bar{X}$

$I$  = Average arrival rate of jobs (general arrival process)

Let  $T$  be the (random) inter-arrival time with (general service time distribution)

$$\begin{cases} E\{T\} = 1/I \\ s_T^2 = E\{T^2\} - [E\{T\}]^2 \end{cases}$$

Let  $X$  be the (random) service time with

$$\begin{cases} E\{X\} = \bar{X} \\ s_X^2 = E\{X^2\} - [E\{X\}]^2 \end{cases}$$

$r = I\bar{X}$  = Traffic Offered

$r < 1$  for queue to be stable

If the mean and variance (or second moment) of the inter-arrival times and the service times are known, then the following bounds have been shown to hold for  $W_q$ , the waiting time in queue, of any G/G/1 queue.

$$\frac{I\mathbf{s}_X^2 - \bar{X}(2 - \mathbf{r})}{2(1 - \mathbf{r})} \leq W_q \leq \frac{I(\mathbf{s}_X^2 + \mathbf{s}_T^2)}{2(1 - \mathbf{r})} \quad (1)$$

*Lower Bound*                      *Upper Bound*

The upper bound of (1) is quite useful. The lower bound is actually not very useful as it often gives a negative result which is a trivial conclusion.

Another interesting (and very useful) bound for the G/G/1 queue has been given for the special case where the inter-arrival time  $T$  satisfies the following property for all values of  $t$ .

$$E\{T - t | T > t\} \leq \frac{1}{I} \quad \text{for all } t \geq 0 \quad (2)$$

- Note that  $E\{T\}=1/I$ . The condition of (2) is not very hard to satisfy. If the inter-arrival time is known to be more than  $t$ , then (2) requires that the expected length of the remaining inter-arrival time should be less than or equal to the *unconditioned* expected inter-arrival time  $1/I$ .
- For the special case when the arrival process is Poisson, the inter-arrival times will be exponentially distributed and will satisfy (2) as an equality.
- Note that many distributions (like say the uniform distribution) will satisfy (2). An exception to this are *hyper-exponential* type distributions

If the arrival process is such that (2) is satisfied, then the following bounds have been shown to hold

$$W_{qU} - \frac{1+r}{2I} \leq W_q \leq W_{qU} \quad (3)$$

where  $W_{qU}$  is the upper bound of (1), i.e.  $W_{qU} = \frac{I(s_X^2 + s_T^2)}{2(1-r)}$



To see the tightness of the bounds of (3), consider using it to find the bounds on  $N_q$ , the mean number waiting in queue for a G/G/1 system.

We get

$$IW_{qU} - \frac{1+r}{2} \leq N_q \leq IW_{qU} \quad (4)$$

- The difference between the upper and the lower bounds is only  $0.5(1+r)$ .
- Note that  $r$  will be small anyway as  $0 < r < 1$  for a stable queue.
- In any case, the difference between the upper and the lower bounds will be between 0.5 and 1. In percentage terms, as  $r \ll 1$ , i.e the traffic increases, this will get increasingly smaller compared to  $N_q$ .

Copyright 2002, Sanjay K. Bose

17

### Heavy Traffic Approximation for the G/G/1 Queue

As  $r \ll 1$ , (i.e. when the offered traffic is high), the distribution of the waiting time in a G/G/1 queue will be approximately an exponentially distributed random variable with mean given by

$$W_q = \frac{I(s_X^2 + s_T^2)}{2(1-r)}$$

Note that this result is an interesting one as it not only provides a mean but also a distribution for the waiting time in queue under very general conditions.

Copyright 2002, Sanjay K. Bose

18

### The G/G/m Queue

$m$  = Number of Servers  
 $r = I\bar{X}$  = Offered Traffic

$W_{q1}$  = Average waiting time in queue  
 for the equivalent G/G/1 queue.

Other notations same as for  
 the G/G/1 queue

#### The Equivalent G/G/1 Queue

Same arrival process of jobs as for the G/G/m queue

For the service times, use  $\frac{\bar{X}}{m}$  and  $\frac{s_x^2}{m^2}$  as the mean and variance,  
 respectively.

Note that the server here works  $m$  times faster than a server in the  
 original G/G/m queue

Copyright 2002, Sanjay K. Bose

19

The following bounds then hold for the G/G/m queue

$$W_{q1} - \frac{(m-1)\bar{X}^2}{2m\bar{X}} \leq W_q \leq I \frac{\left[ s_T^2 + \frac{s_x^2}{m} + \frac{(m-1)(\bar{X})^2}{m^2} \right]}{2 \left[ 1 - \frac{r}{m} \right]}$$

- The value of  $W_{q1}$  may be computed as a lower bound on  $W_q$
- These bounds are rather loose and may not be very useful in practice

Copyright 2002, Sanjay K. Bose

20

### Heavy Traffic Approximation for the G/G/m Queue

For  $(r/m)@1$  in a G/G/m queue, a *heavy traffic approximation result* holds in a manner similar to that given for the G/G/1 case.

Specifically, for  $(r/m)@1$  in a G/G/m queue, the waiting time in queue at steady-state tends towards a random variable with an exponential distribution which has a mean given by -

$$W_q \approx I \frac{\left[ s_T^2 + \frac{s_X^2}{m} \right]}{2 \left[ 1 - \frac{r}{m} \right]}$$

Copyright 2002, Sanjay K. Bose

21

### The M/G/m/m Queue

- Even though the M/G/m queue is hard to analyze, the finite capacity M/G/m/m queue (*m server queue without waiting positions*) is surprisingly easy to analyze. (See additional notes)
- Even more remarkably, its state probability distribution and blocking probability results are identical to those obtained for the corresponding M/M/m/m queue

$$p_k = \frac{\frac{r^k}{k!}}{\sum_{n=0}^m \frac{r^n}{n!}} \quad P_B = B(m, r) = \frac{\frac{r^m}{m!}}{\sum_{j=0}^m \frac{r^j}{j!}} \quad r = I\bar{X}$$

for  $k=0, 1, \dots, m$

Copyright 2002, Sanjay K. Bose

22