1. Introduction

A "network of queues" representation of a multiple-resource model plays an important role in performance analyses of computer/computer-communication systems. Jackson [1] considered a network of m exponential servers labelled as servers 0,1,2,...m-1, in which (i) the job arrival process is Poisson with rate $\lambda(N)$ where N is the total number of jobs found in the system at a given time, (ii) the job routing is characterized by a 1st-order Markov chain, i.e., a job completing service at server i will go next to the service station j with probability $p_{ij}$ irrespective of how this job has reached server i, (iii) the service discipline is any of the so-called "work-conserving" queue disciplines [2], which includes a class of preemptive-resume scheduling disciplines, as well as FCFS and non-preemptive scheduling disciplines, and (iv) the processing rate, $C_j(n_j)$ [work units/sec] of the station j may be an arbitrary function of its local queue size $n_j$. With these assumptions Jackson [1] showed that the equilibrium state distribution of the queue-size vector $\vec{n} = [n_0, n_1, \ldots, n_{m-1}]$ is given in terms of the <u>product</u> of the marginal distributions of the variables $n_j$, j=0,1,2,...,m-1. Namely, the linear difference equation (or the system's balance equation) can be solved via the "separation of variables" method, insofar as the steady-state queue distribution is concerned. Chandy [3] calls the linear difference equations for the individual coordinate variables "local balance" equations.

Recently, Chandy [3], Baskett and Muntz [4] and Baskett et al. [5] have made a substantial extension of the Jackson's model. Their result can be summarized as follows: if the assumption (iii) is replaced by a more strict one, namely if the service-discipline of a given station is either the (round-robin) processor sharing (PS) or preemptive-resume "last-come, first-served" (LCFS) discipline, then

ON GENERALIZATION OF JOB ROUTING BEHAVIOR IN A QUEUEING
NETWORK MODEL

Hisashi Kobayashi and Martin Reiser

IBM Thomas J. Watson Research Center
Yorktown Heights, New York

Typed by Mrs. Sue Hritz on IBM MT/ST.

ABSTRACT: Prior work on queueing network models has always
assumed that job routing behavior is governed by a 1st-order
Markov chain. In the present paper we eliminate this re-
striction and allow any type of probabilistic routing that
is representable in terms of a Markov chain of arbitrary
order. We will show that the simple product form solution
obtained earlier for more restrictive models remains to hold
here. We then obtain an exceedingly simple result: the
only parameter that appears in the queue-size distribution
is a set of values $\{W_j\}$, where $W_j$ is the expected workload
that a job places on server j during the job's entire life.

the service-time distribution at this station can be a general distribution. This rather surprising result is due to the fact that with the PS or LCFS discipline, an M/G/1 queue is essentially converted into an "equivalent" M/M/1 system, insofar as the departing process and the steady state queue distribution are concerned [6]. That is, the queue distribution is geometric and job departures exhibit a Poisson process*. They have also shown that if the queue-dependent service rate of the assumption (iv) is specialized as $C_j(n_j) = n_j \cdot C_j^*$, then the service-time distribution of station j can be also general. This is because the service station with $C(n) = n \cdot C^*$ is equivalent to infinitely many servers (IS) each of which has the processing rate C* [work units/sec]. It is also known that an M/G/∞ queue exhibits a Poisson departure. The queue size distribution in this case is not a geometric distribution but a Poisson distribution. The above authors [4-6] have also showed that those servers with the PS or LCFS discipline or IS stations allow multi-class jobs in terms of both the service time distribution and job routing transitions.

Kobayashi and Reiser [7] subsequently treated the case in which the underlying Markov chain is decomposable to multiple subchains. Furthermore, they have developed an efficient computational algorithm [8] which determines the normalization constant of such analytical solutions.

All these works referenced above, however, maintain the assumption (ii), i.e., the job transition behavior is governed by a 1st-order Markov chain. In the present paper we will essentially remove the assumption (ii) and will show that the class of queueing network models which satisfy the "separation of variables" can be substantially enlarged. This result will bring significant implications into our system modeling practice.

---

* To be more precise, the server adopts the PS or LCFS discipline, but its processing rate $C_j(n_j)$ is independent of $n_j$. If the queue-dependent processing rate is introduced, the departure process becomes a state-dependent Poisson process, and the queue is no longer geometric.

II. A Cyclic Queue with a General Distribution for the Number of Cycles.

Let us start with the simplest queueing network, namely, a two-stage cyclic queue as shown in Figure 1. We define a "cycle" to be a routing of a job from the branching point A to server 1, server 0, and back to A. Then the number of cycles k which a given job makes during its lifetime in the system is

$$p_k = (1-\alpha)\alpha^k, \quad k = 0,1,2,\ldots \tag{1}$$

where $\alpha$ is the probability that a job cycles back to server 1 from the point A. Thus the assumption of 1st-order Markov transitions in a cyclic queue implies that the variable k is geometrically distributed. The model of Figure 1 is often used in the analysis of a multiprogrammed computer system, even though the distribution of (1) is not always justifiable.

Consider now a discrete distribution $\{p_k, k \geq 0\}$ of general form, an example of which is shown in Figure 2. We define the probability generating function (p.g.f.) $P(z)$ by

$$P(z) = \sum_{k=0}^{\infty} p_k z^k \tag{2}$$

Practically speaking, for any given $\{P_k\}$, the p.g.f. $P(z)$ can be written as a rational function of z:

$$P(z) = \frac{R(z)}{Q(z)} \tag{3}$$

Then we can obtain the following expansion:

$$P(z) = b_0 + \sum_{r=1}^{q} a_0 a_1 \cdots a_{r-1} \cdot b_r \prod_{i=1}^{r} \frac{1-\alpha_i}{1-\alpha_i z} \tag{4}$$

where q is the degree of the polynomial $Q(z)$, and $\{\alpha_i^{-1}\}$ are the characteristic roots of $Q(z) = 0$. The coefficients $\{a_i\}$ and $\{b_i\}$ satisfy

$$a_i + b_i = 1, \quad 0 \le i \le q-1 \tag{5}$$

and

$$b_q = 1 \tag{6}$$

The representation (4) can be schematically shown in Figure 3, which is equivalent to cascaded geometric distributions with parameters $\alpha_1, \alpha_2, \dots, \alpha_q$.

The expansion (4) is a discrete analogue of the exponential stage representation [12] of a general service time whose Laplace transform is a rational function of the Laplacian variable s. In general, this representation involves the formal use of complex transition probabilities, since the characteristic roots $\alpha_i^{-1}$ of $Q(z)$ can take on complex values.

Examples

(1) Consider the distribution

$$P_k = \binom{k + q - 1}{q-1} (1-\alpha)^q \alpha^k, \qquad k=0,1,2,\dots$$

which is a shifted version of the negative binomial or Pascal distribution. Its p.g.f. is then

$$P(z) = \left( \frac{1-\alpha}{1-\alpha z} \right)^q$$

Thus the expansion (4) is simply obtained as $a_0 = a_1 = \cdots = a_{q-1} = 1$, $b_0 = b_1 = \cdots = b_{q-1}$, $b_q = 1$ and $\alpha_1 = \alpha_2 = \cdots = \alpha_q = \alpha$.

(2) A Poisson distribution

$$P_k = \frac{\lambda^k}{k!} e^{-\lambda}, \qquad n \ge 0$$

has the corresponding p.g.f.

$$P(z) = e^{\lambda(z-1)}$$

which is clearly not a rational function of z. If we allow, however, the formal passages to the limit we can represent, P(z) as the infinite stage:

$$\lim_{q \to \infty} \left[ \frac{1 - \frac{\lambda}{k}}{1 - \frac{\lambda z}{q}} \right]^q = e^{\lambda(z-1)}$$

(3) Let $\{p_k\}$ be of oscillatory form with geometric decay as illustrated in Figure 2.

$$p_k = C \cdot r^k (1 - \cos\theta n)$$

where $0 < r < 1$ and C is the normalization constant such that $\sum_k P_k = 1$. The p.g.f. is

$$P(z) = \frac{C(1-\cos\theta)(1+rz)rz}{(1-rz)(1+r^2z^2 - 2rz \cos\theta)}$$

Thus one of the probability parameters $\alpha_1 = r$ is real, but the other two are complex conjugate $\alpha_2 = re^{i\theta}$, $\alpha_3 = re^{-i\theta}$. It can be shown that although the probabilities associated with the individual fictitious stages may be complex, the probabilities associated with real states are real.

In Figure 3 each box labelled z corresponds to one cycle in Figure 1. Corresponding to the q fictitious stages defined above, we now introduce q different classes. A job whose routing is in its rth stage is classified as a class-r job, $1 \leq r \leq q$. For notational convenience we define that all entering jobs proceed first to server 0 as class-0 jobs. A job, after receiving its first service at server 0, either leaves the system immediately, or becomes a class-1 job and cycles around the servers 1 and 0 as many as $k_1$ times, where

the random variable $k_1$ is geometrically distributed:

$$p_{k_1} = (1-\alpha_1)\alpha_1^{k_1} \qquad k_1 = 0,1,2,\ldots \tag{7}$$

The job then leaves the system, otherwise changes its status to a class-2 job and cycles around $k_2$ times, and so forth. Thus the average number of visits a job makes to servers j with class membership r during its entire life in the system is

$$e_{j,r} = a_0 a_1 \cdots a_{r-1} \frac{\alpha_r}{1-\alpha_r} \quad , \quad j=0,1 \quad 1 \le r \le q \tag{8}$$

For r=0, we have:

$$e_{0,0} = 1 \quad \text{and} \quad e_{1,0} = 0 \tag{9}$$

Let us define the system state vector by

$$\vec{n} = [\underline{n}_0, \underline{n}_1] \tag{10}$$

where

$$\underline{n}_j = [n_{j,0}, n_{j,2}, \cdots, n_{j,q}], \quad j=0,1 \tag{11}$$

with

$$n_{1,0} = 0 \tag{12}$$

If the job arrival process from the outside is a Poisson process with rate $\lambda(N)$, where $N = |\vec{n}|$ is the total number of jobs in the system, then the equilibrium state distribution is given [5,7,8] by

$$p(\vec{n}) \begin{cases} = c \, \Lambda(|\vec{n}|) g_0(\underline{n}_0) g_1(\underline{n}_1), & \text{if } \vec{n} \text{ is feasible,} \\ = 0, & \text{otherwise} \end{cases} \tag{13}$$

where c is the normalization constant, and $\Lambda(N)$ is determined solely by the arrival rate function $\lambda(n)$:

$$\Lambda(N) = \prod_{n=0}^{N-1} \lambda(n) \tag{14}$$

If the arrival rate does not depend on the number of jobs in the system, i.e., $\lambda(n)=\lambda$, then clearly $\Lambda(N)=\lambda^N$. If the network is a closed network, then we define $\Lambda(N)=1$. The function $g_j(\underline{n}_j)$ is the (improper) probability distribution of $\underline{n}_j$ and takes the form [8]

$$g_j(\underline{n}_j) = D_j(n_j) \, n_j! \, \prod_{r=0}^{q} \left[ \frac{1}{n_{jr}!} \, (e_{jr}\overline{w}_j)^{n_{jr}} \right], \quad j=0,1 \tag{15}$$

if server j uses the round-robin processor sharing (PS) scheduling or the pre-emptive-resume "last-come, first-served" (LCFS) scheduling. The parameter $\overline{w}_j$ of (15) is the average service demand [work units/job] a job places upon server j. The function $D_j(n)$ of (15) is dependent only on the queue-dependent processing rate function $C_j(n)$ [work units/sec] of the server j and is defined by

$$D_j(n_j) = \prod_{n=1}^{n_j} \frac{1}{C_j(n)} \quad . \tag{16}$$

If the processing rate is independent of the queue size, i.e., $C_j(n) = C_j{}^*$. then simply $D_j(n) = 1/C_j{}^{*n}$. Practical examples in which we can make use of the queue-dependent processing rate are as follows. The first example is representation of a multiprocessor: if server j is a symmetric m-way multiprocessor (m-parallel server) with each processing rate $C_j{}^*$ , then we can define $C_j(n)= \min\{n,m\}C_j{}^*$. The infinite server (IS) defined earlier is the limiting case $m \to \infty$. The second example is consideration of the processor degradation due to system overhead. In

general, the effective rate of the processor will decrease as the number of
jobs in the server increases. The PS and LCFS disciplines defined above, for
examples, involves "task switching" and that clearly should introduce some loss
in the processor's productivity and this factor can be suitably accounted for
in terms of $C_j(n)$. It will be clear that an asymmetric multiprocessor, or
multiprocesssor with degradation factor can be suitably represented as $C_j(n) =$
min $\{n,m\}$ $C_j^*(n)$. In the queueing literature [1,4-7] we usually use, instead
of $C_j(n)$, the quantity $\mu_j(n) = C_j(n)/\overline{w}_j$ [jobs/sec] which is also called queue-
dependent service rate. Needless to say, if the service demand on server j is
exponentially distributed with mean $\overline{w}_j$, the formula (15) is valid not only for
PS and LCFS disciplines, but also for any work-conserving queue discipline, in-
cluding FCFS discipline.

If the station j is an IS (infinitely many servers), we can write $C_j(n) =$
$nC_j^*(n)$. Then $g_j(\underline{n}_j)$ is given, by rewriting (15), as

$$g_j(\underline{n}_j) = D_j^*(n_j) \prod_{r=0}^{q} \left[ \frac{1}{n_{jr}!} (e_{jr}\overline{w}_j)^{n_{jr}} \right] \quad , \quad j=0,1 \quad (17)$$

where

$$D_j^*(n_j) = \prod_{n=1}^{n_j} \frac{1}{C_j^*(n)} \quad (18)$$

If the server j is an IS with no degradation, it follows that $D_j^*(n)=1/C_j^{*n}$.
The formula (17) holds for an arbitrary distribution of service demand $w_j$.

The distribution of our interest is the marginal distribution of the
variables.

$$n_j = \sum_{r=0}^{q} n_{jr} \quad (19)$$

for j=0,1. From (14) we readily obtain that

$$p(n_0,n_1) \begin{cases} = c\Lambda(n_0+n_1) \, f_0(n_0)f_1(n_1), & \text{if } (\dot n_0,n_1) \text{ is feasible} \\ = 0, & \text{otherwise,} \end{cases} \qquad (20)$$

For a given value of $n_j$ the function $g_j(\underline{n}_j)$ of (15), excluding the term $D_j(n)$, is proportional to a multinomial distribution. In general, let $\underline{x} = [x_1 x_2 \ldots x_q]$ have a distribution given by

$$g(\underline{x}) = c(x_1+x_2+\ldots+x_q)! \prod_{r=1}^{q} \frac{1}{x_r!} \, \rho_r^{x_r}, \text{ for all } x_r \geq 0 \qquad (21)$$

Then the sum $y = x_1+x_2+\ldots+x_q$ has the distribution

$$f(y) = c \, \rho^y, \quad y \geq 0 \qquad (22)$$

which is a geometric distribution, where $\rho = \sum_r \rho_r$ and the normalization constant is determined as $c = 1 - \rho$. Hence the distribution (20) may be called a multinomial distribution **built on** a geometric distribution. By applying the formulae (20) and (21) to (15) we obtain the marginal distribution of $n_j$.

$$f_j(n_j) = D_j(n_j) \left( \sum_{r=0}^{q} e_{jr} \bar w_j \right)^{n_j} \qquad (23)$$

if server j adopts PS, LCFS, FCFS, etc. If server j is an IS, $g_j(n_j)$ takes the form of multiple Poisson distribution as shown in (17). In general if $\underline{x}$ has the distribution

$$g(\underline{x}) = \prod_{r=1}^{q} \frac{\lambda_r^{x_r}}{x_r!} \, e^{-\lambda_r} \qquad (24)$$

Then the sum $y = \sum_r x_r$ has the Poisson distribution

$$f(y) = \frac{\lambda^y}{y!} \, e^{-\lambda} \qquad (25)$$

where $\lambda = \sum_r \lambda_r$. Then applying the formulas (23) and (24) to (17), we obtain

$$f_i(n_i) = \frac{D_j^*(n_j)}{n_j!} \left( \sum_{r=0}^{q} e_{jr}\bar{w}_j \right)^{n_j} \tag{26}$$

if server j is an IS.

From (8) - (10) we have

$$e_1 = \sum_{r=0}^{q} e_{1r} = \sum_{r=1}^{q} a_0 a_1 \cdots a_{r-1} \cdot \frac{\alpha_r}{1-\alpha_r} \tag{27}$$

which turns out to be $E[k]$, the average number of cycles. This can be shown by using the identity

$$E[k] = P'(1) = \sum_{r=1}^{q} a_0 a_1 \cdots a_{r-1} b_r \cdot \sum_{i=1}^{r} \frac{\alpha_i}{1-\alpha_i}$$

$$= \sum_{i=1}^{q} a_0 a_1 \cdots a_{i-1} \frac{\alpha_i}{1-\alpha_i} \sum_{r=i}^{q} a_i \cdots a_{r-1} b_r$$

$$= \sum_{i=1}^{q} a_0 a_1 \cdots a_{i-1} \frac{\alpha_i}{1-\alpha_i} \tag{28}$$

The last expression was obtained by repetitive use of the identity $a_j + b_j = 1$ for $1 \leq j \leq q-1$ and $b_q = 1$. Therefore from (26) and (27)

$$e_1 = E[k] \tag{29}$$

Similarly we obtain the average number of visits that a job makes to server 0.

$$e_0 = E[k] + 1 \tag{30}$$

Thus, the distribution (23) and (26) can be rewritten as

$$f_i(n_j) = D_j(n_j) W_j^{n_j} \quad \text{for PS, LCFS, FCFS} \tag{31}$$

and

$$f_i(n_i) = \frac{D_j^*(n_j)}{n_j!} W_j^{n_j} \quad \text{for IS} \tag{32}$$

where

$$W_j = e_j \; \overline{w}_j \quad \text{[work units/job]} \qquad \qquad (33)$$

represents the total average workload that a job places on the server j during the total life time of the job.

In the above derivation we assumed that the service demands placed on the server j by a job in its successive cycles are i.i.d. random variables, and hence have the common mean $\overline{w}_j$. By a straightforward extension of the argument, however, we can show that the solution forms (31) and (32) hold for PS, LCFS and IS even when the service demands are not identically distributed. What counts in the queue size distribution is the total average work a job brings in, and how the total work is distributed in individual cycles is immaterial. This is a surprisingly strong result. We have to stress again, however, that in order for the formulae (31), (32) to hold for FCFS or any work-conserving discipline, the service demands in successive cycles must be drawn from identical exponential distributions.

## III.  Job Routing Characterized by a High-Order Markov Chain

The notion of "cycles" has some difficulty when we attempt to extend it to a queueing network with general topology. We will therefore take a different approach in this section:  the first-order Markov chain, which characterizes Jackson's model [1] and other related work [3-8], will now be replaced by a high-order Markov chain. Let us start again with a simple network with two servers which we now again denote by server 0 and server 1. In the 1st-order Markov model, the transition probabilities $p_{s,s'}$ were defined over s,s'=0,1. Now let us assume that job routing is statistically characterized by a 2nd-order Markov chain. Then the probabilities $p_{s,s'}$ are now defined over states s,s'=(00), (01), (10), (11).

A job is said to be in state (ij) if the job is now at server j just after completing service at server i. For notational conciseness we use integers 0, 1, 2 and 3 to denote the states (00) (01) (10) and (11), respectively. Therefore, jobs in either states 0 or 2 are located at server 0 and those in states 1 or 3 are at server 1. By using the discrete time (or step) parameter k, a job routing can conveniently be represented by a "trellis" picture of Figure 4.

In Figure 4 we introduce an additional state, s=4, which is an absorbing state. A transition to state 4 at step k means that the job leaves the system after k services. For example, the path $0 \to 1 \to 3 \to 2 \to 4$ in Figure 4 means that a job enters server 0 first, moves to server 1, and again to server 1, and finally goes back to server 0 and then leaves the system. Let us denote by $e_s(k)$ the probability that a job is in state s at step k, $\sum\limits_{s=0}^{4} = 1$ for all k. Then the equilibrium state distribution of the system state vector

$$\vec{n} = \{n_{s,k}; \quad 0 \leq s \leq 3; \quad k = 0, 1, 2,...\} \tag{34}$$

is given again by eq. (14), where $g_j(\underline{n}_j)$, j=0, 1 should now be interpreted as the (improper) distribution of subvectors,

$$\underline{n}_0 = \{n_{sk}; \quad s=0,2; \quad k=0,1,2,...\} \tag{35}$$

and

$$\underline{n}_1 = \{n_{sk}; \quad s=1,3; \quad k=0,1,2,...\} \tag{36}$$

respectively. Clearly $\underline{n}_0$ and $\underline{n}_1$ together consitute the system state vector $\vec{n}$. The function $g_i(\underline{n}_j)$ takes essentially the same form as (16) or (18):

$$g_0(\underline{n}_0) = D_0(n_0) \, n_0! \prod_{s=0,2} \prod_{k=0}^{\infty} \left[ \frac{1}{n_{sk}!} \, (e_s(k)\bar{w}_0)^{n_{sk}} \right]$$

$$\text{for PS, LCFS, FCFS, etc.} \tag{37}$$

and

$$g_0(\underline{n}_0) = D_0^*(n_0) \prod_{s=0,2} \prod_{k=0}^{\infty} \left[ \frac{1}{n_{sk}!} (e_s(k)\bar{w}_0)^{n_{sk}} \right]$$

for IS. (38)

The distribution $g_1(\underline{n}_1)$ is the same as $g_0(\underline{n}_0)$ except that the product is taken over $s=1,3$ instead of $s=0,2$.

The distributions of $n_0$ and $n_1$ are then derived as marginal distributions of $p(\vec{n})$, resulting in the same form as (20) where $f_0(n_0)$ is readily obtained from (37) as

$$f_0(n_0) = D_0(n_0) \left( \sum_{s=0,2} \sum_{k=0}^{\infty} e_s(k) \bar{w}_0 \right)^{n_0}.$$

for PS, LCFS, FCFS, etc. (39)

or from (38) as

$$f_0(n_0) = \frac{D_0^*(n_0)}{n_0!} \left( \sum_{s=0,2} \sum_{k=0}^{\infty} e_s(k)\bar{w}_0 \right)^{n_0}$$

for IS. (40)

The quantity $e_0$ defined by

$$e_0 = \sum_{s=0,2} \sum_{k=0}^{\infty} e_s(k) \tag{41}$$

is the total average number of visits a job makes to server 0 during the job's life time. Similarly, a job visits server 1, on the average, as many as

$$e_1 = \sum_{s=1,3} \sum_{k=0}^{\infty} e_s(k) \tag{42}$$

times. By using $W_j$ defined in (33), we are led again to the simple expressions (31) and (32).

As we discussed in the previous section, we can further generalize the definition (33). If the server j uses the PS or LCFS discipline or server j is an IS, then we can allow that the service demand distributions be different at different states s and at different steps k. Letting $\overline{w}_s(k)$ be the average service requested by a job in state s at time k, we now define

$$W_0 = \sum_{s=0,2} \sum_{k=0}^{\infty} e_s(k)\overline{w}_s(k) \tag{43}$$

and

$$W_1 = \sum_{s=1,3} \sum_{k=0}^{\infty} e_s(k)\overline{w}_s(k) \tag{44}$$

The restrictions associated with FCFS or arbitrary work-conserving discipline are the same as discussed at the end of the previous section.

The evaluation of $e_j$, j=0,1 can be done by a straightforward application of the Markov chain theory. The probabilities $\{e_s(k)\}$ satisfy the equation

$$e_s(k) = \sum_{s'=0}^{4} e_s(k-1)p_{s's} \tag{45}$$

for $0 \le s \le 4$ and $k \ge 1$. By defining a row vector

$$\underline{e}(k) = [e_0(k),\ e_1(k),\ldots,e_4(k)], \tag{46}$$

and the corresponding generating function vector

$$\underline{E}(z) = \sum_{k=0}^{\infty} \underline{e}(k)z^k \tag{47}$$

we obtain from (45)

$$\underline{E}(z) - \underline{e}(0) = z\ \underline{E}(z)\ \underline{P} \tag{48}$$

where $\underline{P}$ is the matrix $[P_{s's}]$. We can then derive the well-known formula

$$\underline{E}(z) = \underline{e}(0) \ [\underline{I} - z\underline{P}]^{-1} \qquad\qquad (49)$$

Denoting the individual components of $\underline{E}(z)$ as $E_s(z)$, $0 \leq s \leq 4$, we have

$$e_0 = E_0(1) + E_2(1) \qquad\qquad (50)$$

and

$$e_1 = E_1(1) + E_3(1) \qquad\qquad (51)$$

Example:

Consider the simple cyclic queue of Figure 1. Because of this special configuration, states $0 = (00)$ and $3 = (11)$ are never reached. So it suffices to consider only the three states, $s = 1, 2$ and $4$. By defining

$$\underline{E}(z) = [E_1(z), E_2(z), E_4(z)]$$

and

$$\underline{P} = \begin{bmatrix} 0 & 1 & 0 \\ \alpha & 0 & 1-\alpha \\ 0 & 0 & 1 \end{bmatrix}$$

we obtain

$$\underline{E}(z) = \underline{e}(0) \begin{bmatrix} 1 & -z & 0 \\ -\alpha z & 1 & -(1-\alpha)z \\ 0 & 0 & 1-z \end{bmatrix}^{-1}$$

Since all jobs enter from server 0, the initial condition is

$$\underline{e}(0) = [\,0 \quad 1 \quad 0\,]$$

which leads to

$$E_1(z) = \frac{\alpha z}{1-\alpha z^2} \ ,$$

and

$$E_2(z) = \frac{1}{1-\alpha z^2}$$

Hence

$$e_0 = E_2(1) = \frac{1}{1-\alpha}$$

and

$$e_1 = E_1(1) = \frac{\alpha}{1-\alpha}$$

## IV. Extensions to a General Network Topology with Markov Chain of Higher Order

The presentation of Section III was made by choosing the simplest example, i.e., a network of two servers and the 2nd-order Markov chain. Its extension to a queueing network with general topology with job routings characterized by a higher-order Markov chain is now straightforward. If there are $m$ service centers $0,1,2,\ldots,m-1$ in a network and the job routing transitions are characterized by an hth-order Markov chain, there are $m^h$ different states a job can take on, which we denote as before by integers $s = 0,1,2,\ldots m^h-1$. We then form a trellis picture with $m^h$ different states plus an absorbing state which we denote by $s=m^h$. We define $e_s(k)$ as before for $s=0,1,2,\ldots,m^h$ and $k=0,1,2,\ldots$ We then define parameters

$$e_j = \sum_{\substack{s=j \\ (\bmod\ m)}} \sum_{k=0}^{\infty} e_s(k) \tag{52}$$

where the summation over s is taken over those s which satisfy

$$s(\text{modulo } m) = j \tag{53}$$

For example $e_0$ is obtained by summing over $s=0,m,2m,\ldots,m^{h-1}$. Similarly, $e_1$ is the sum of the terms with $s=1,m+1,2m+1,\ldots,m^{h-1}+1$, and so forth. By defining

the $m^h$-dimensional p.g.f. $\underline{E}(z)$ as in (47) and using the formula (49), we can evaluate $\dot{e}_j$ by

$$e_j = \sum_{\substack{s=j \\ (\text{mod } m)}} E_s(1) , \quad j = 0, 1, \ldots , m-1 \tag{54}$$

where $E_s(z)$'s are the elements of $\underline{E}(z)$. The entire results of Section III carry over to this general case in an obvious way, and thus the joint queue size distribution of this general network is given by

$$p(n_0\ n_1, \ldots, n_{m-1}) = c\Lambda(N) \prod_{j=0}^{m-1} f_j(n_j) \tag{55}$$

where

$$N = \sum_{j=0}^{m-1} n_j \tag{56}$$

and

$$f_j(n_j) = D_j(n_j)\ W_i^{n_j} \quad \text{when server } j \text{ adopts PS, LCFS, FCFS, etc.} \tag{57}$$

and

$$f_j(n_j) = \frac{D_j^*(n_j)}{n_j!}\ W_j^{n_j} \quad , \text{ when server } j \text{ is an IS.} \tag{58}$$

for $j = 0, 1, 2, \ldots , m-1$.

Throughout the above discussion we assumed that the network is open, and thus jobs arrive with rate $\lambda(n)$. It is not difficult to extend the above result to a closed network as was done in the earlier work [5,7,8]. As mentioned earlier, we define in this case $\Lambda(n)=1$. The parameters $\{e_j\}$ are not determinable up to a common scaling factor. If we choose the factor such that $e_{j*}=1$ for some $j*$, then $\{e_j\}$ represents the average number of visits that a job makes to server $j$ between its consecutive visits to server $j*$. The workload parameter $W_j$ then represents the expected total amount of work that the job brings into server $j$ during that cycle.

Throughout our discussion thus far we assumed a homogeneous job population in the sense that the routing behavior of each job is characterized by the same transition $p_{s,s'}$. The formulue (57) and (58) remain to hold even if multiple classes r=1, 2, ..., R, of jobs are introduced, as long as their arrival processes are mutually independent Poisson processes with rate $\lambda_r$, r=1,2,...,R. The introduction of R classes, is essentially equivalent to extending the dimension of state space by factor of R. Multiple Poisson streams can be created from a single Poisson stream of rate $\lambda = \sum_{r=1}^{r} \lambda_r$, which branches out to one of R parallel streams with probabilities $\{\lambda_r/\lambda\}$.

The simple results of (57) and (58) are essentially dependent on one important condition. Namely, an arrival process to each state (not server) must be Poisson. Here the distinction between state and server is important. In the cyclic queue of Figure 1, job arrival process to server 0, for example, is not Poisson. This is because some of the arrivals to server 0 are due to past departures from the server 0 and hence the present arrival is dependent on past arrivals. If we distinguish the states of job in its different cycles, the arrival of a job to a given state is Poisson. If the condition of Poisson arrivals is met over a suitably defined state space, then the product form solution holds and its marginal distribution is reduced to the simple formulae (54) - (57). Thus a crucial step is to examine whether we can define a state space in which the Poisson arrival condition is satisfied at each state.

## IV.  Conclusions

In this paper we have shown that a simple product form solution of a queueing network model can be obtained, even when the 1st-order Markov assumption on job routing is removed.  This result significantly enlarges a class of problems to which we can realistically apply a "network of queue" model. We have shown that for any job routing behavior (which can be represented in terms of a Markov chain of arbitrary order) the only parameter that appears in the queue-size distribution is $W_j$, the average workload that a job places on server j.  This result also simplifies our efforts on model constructions and validations which use empirical data; namely, we need not measure detailed transition probabilities, but simply estimate the total workload requested of each resource by jobs.

References

[1]  J. R. Jackson, "Job Shop-Like Queueing Systems", Management Science, Vol. 10, 1963, pp. 131-142.

[2]  L. Kleinrock, "A Conservation Law for a Wide Class of Queueing Disciplines", Naval Res. Log. Quart. Vol. 12, pp. 181-192, 1965.

[3]  K. M. Chandy, "The Analysis and Solution for General Queueing Networks", Proceedings of the Sixth Annual Princeton Conference on Information Sciences and Systems.  (Princeton University, N.J.) March 1972.

[4]  F. Baskett and R. R. Muntz, "Queueing Network Models with Different Classes of Customers", Proceedings of the 6th Annual IEEE Computer Society International Conference, September 1972, pp. 205-209.

[5]  F. Baskett, K. M. Chandy, R. R. Muntz, and F. G. Palacios, "Open, Closed, and Mixed |Networks of Queues with Different Classes of Customers", to appear in Journal of the ACM.

[6]  R. R. Muntz, "Poisson Departure Processes and Queueing Networks", IBM Research Report RC-4145, September 1972.  Also in Proceedings of the Seventh Annual Princeton Conference on Information Sciences and Systems, March 1973, pp. 428-434.

[7]  H. Kobayashi and M. Reiser, "Some Results on Queueing Models with Different Classes of Customers", Proceedings of the Eighth Annual Princeton Conference on Information Sciences and Systems, March 1974.

[8]  M. Reiser and H. Kobayashi, "Queueing Networks with Multiple Closed Chains: Theory and Computational Algorithms", IBM Research Report RC 4919, July, 1974, IBM T. J. Watson Research Center, Yorktown Heights, New York.  Also to appear in IBM Journal of Research and Development, 1975

[9]  D. R. Cox, "A Use of Complex Probabilities in the Theory of Stochastic Processes", Proc. Camb. Phil. Soc., Vol. 51, 1955, pp. 313-319.
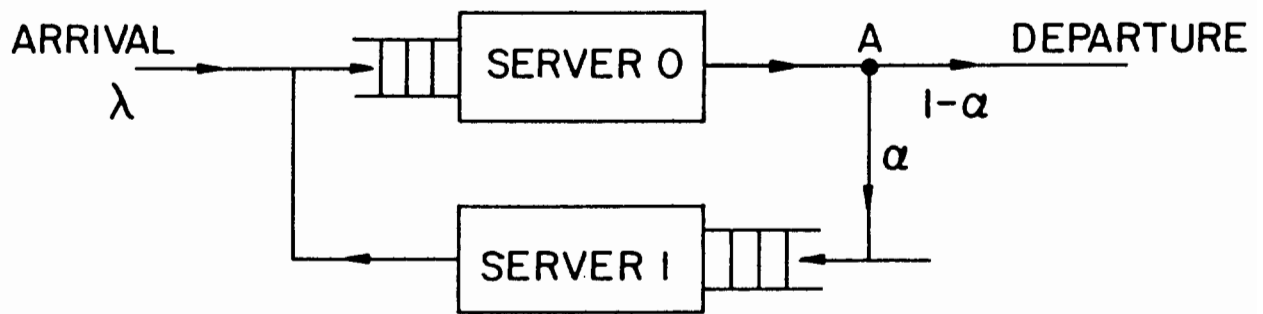
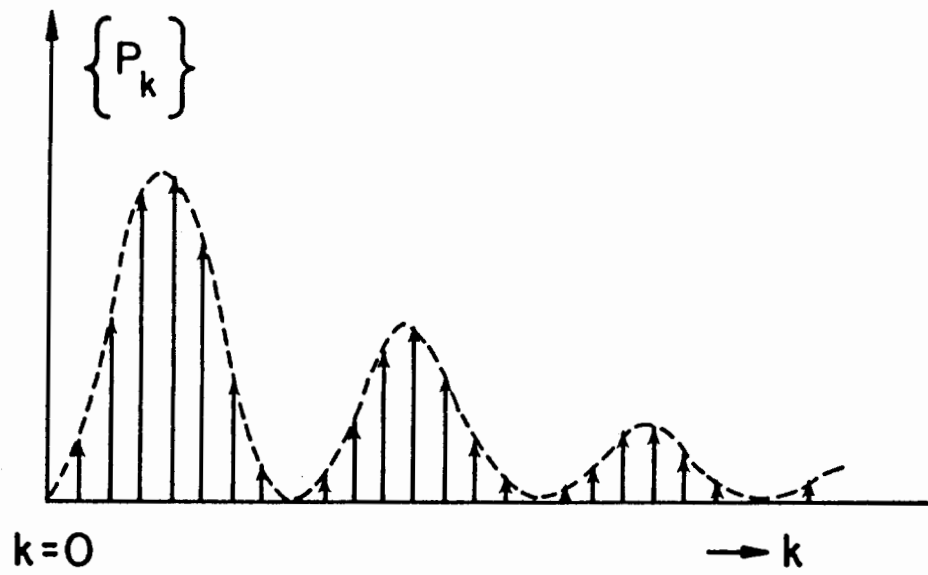Figure 1. A Cyclic Queueing System.

Figure 2.   An Example of General Distribution of the Number
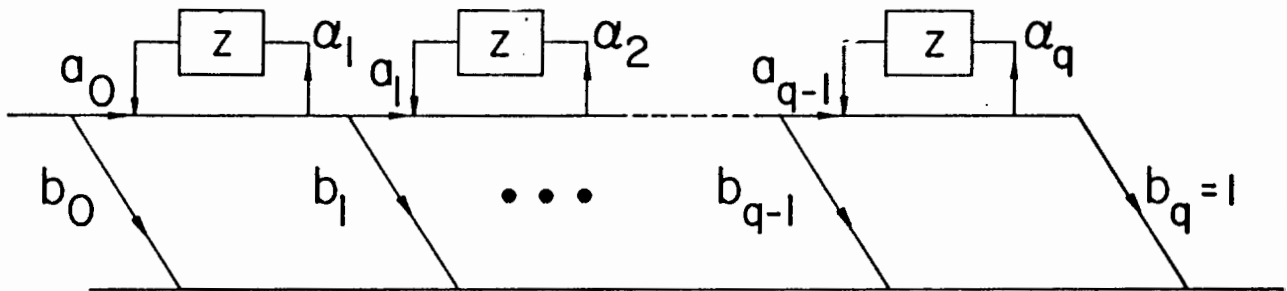            of Cycles, $\{p_k\}$.

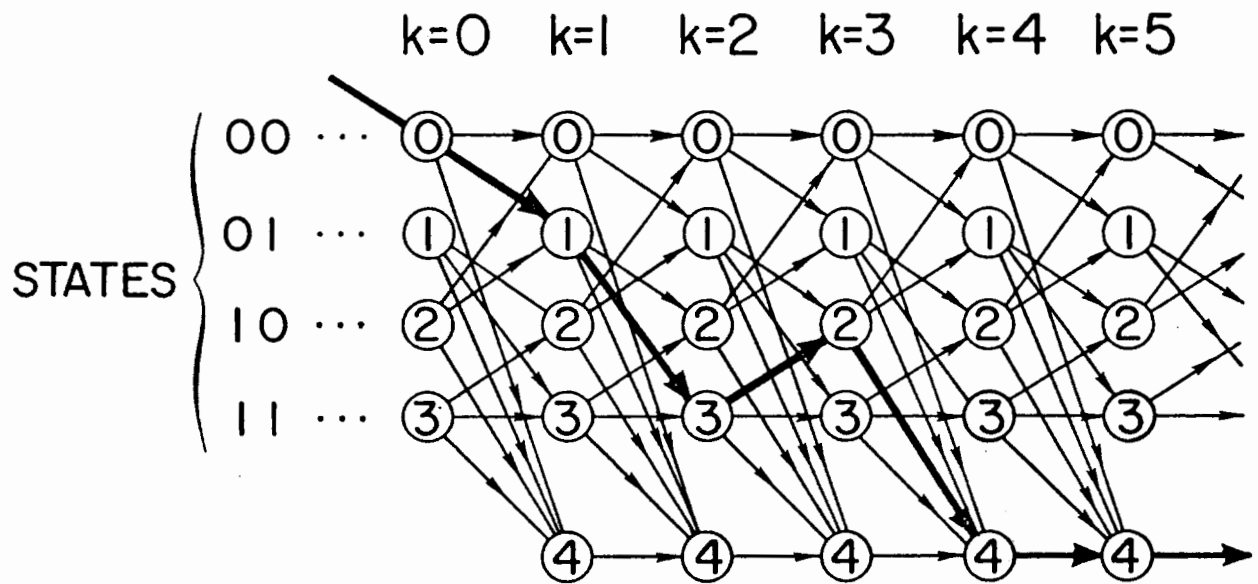Figure 3. The Schematic representation of Cascaded Geometric Distribution of P(z).

Figure 4. Trellis Picture Representation of Job Routing Transition.