# Open Networks of Queues: Their Algebraic Structure and Estimating Their Transient Behavior

William A. Massey

*Advances in Applied Probability* is currently published by Applied Probability Trust.

# OPEN NETWORKS OF QUEUES: THEIR ALGEBRAIC STRUCTURE AND ESTIMATING THEIR TRANSIENT BEHAVIOR

WILLIAM A. MASSEY,* *Bell Laboratories*

### Abstract

We develop the mathematical machinery in this paper to construct a very general class of Markovian network queueing models. Each node has a heterogeneous class of customers arriving at their own Poisson rate, ultimately to receive their own exponential service requirements. We add to this a very general type of service discipline as well as class (node) switching. These modifications allow us to model in the limit, service with a general distribution. As special cases for this model, we have the product-form networks formulated by Kelly, as well as networks with priority scheduling. For the former, we give an algebraic proof of Kelly's results for product-form networks. This is an approach that motivates the form of the solution, and justifies the various needs of local and partial balance conditions.

For *any* network that belongs to this general model, we use the operator representation to prove stochastic dominance results. In this way, we can take the transient behavior for very complicated networks and bound its joint queue-length distribution by that for *M/M/*1 queues.

FOCK SPACES; FREE SEMIGROUPS; STOCHASTIC DOMINANCE; KELLY NETWORKS; TENSOR GEOMETRIC DISTRIBUTION

## 1. Introduction

In this paper we develop a complete operator theory for discrete-state Markovian queueing networks. These queueing systems are governed by an associated linear operator via the Kolmogorov forward equations for general Markov processes. We construct the operator for the particular queueing network in a manner that exploits the intrinsic algebraic structure of these queueing models. This in turn gives us a more 'hands on', direct approach for dealing with these systems, as we shall demonstrate.

The *M/M/*1 queue and the Jackson network are the basic corner-stones of queueing theory for single-server systems and queueing networks respectively. We shall extend these models to incorporate the complexity of different classes of customers. For a single-server queue, this means that each class has its own

Poisson arrival rate and its own exponential service distribution. The different requirements of each class force us to be aware of the order in which customers are served. This is something that we never do for the $M/M/1$ queue when we consider its queue-length distribution. We need to define a general discipline that determines where an arriving customer is inserted in line. It should also determine which served customer is removed. Our description of the discipline will be general enough to be far more than adequate to describe first-come–first-served (FIFO), last-come–first-served (LIFO), processor sharing (PS), and even priority scheduling.

In addition, we allow for a special type of feedback which we shall call *internal switching*. A customer having completed service will then hold its position in line and change its class, with some given probability. These class changes are governed by a probability switching matrix, which has the customer eventually leaving the queue with probability 1. In this manner we can model, for example, customers whose service time has a phase-type distribution.

The general network that we construct is a collection of single-server queues as described above. Each queue for each node has its own type of discipline. We merely add to this the notion of *external switching*. In addition to the usual switching from one node to another as in the Jackson network, we allow for the class of the customer to change also. Unlike internal switching, the customer, when sent to the new node, is subjected to the insertion rule for that node.

In Section 2 we briefly review the functional analytic tools that we use to represent the probability distributions of various queueing models. After that, we undertake a fourfold plan to build up to our general network model. In Section 3, we review the typical construction of the $M/M/1$ queueing system. We then reinterpret it in operator notation. In so doing, we introduce the fundamental notions of right- and left-shift operators. We also introduce the two basic themes of the paper. First, we take the local and partial balance conditions that occur for product-form networks, and reinterpret them in simple algebraic terms. This leads to a purely algebraic proof for the product-form networks discussed in Kelly [2]. The method of proof is one that is elegant and motivates the various hypotheses for these systems. Second, we develop stochastic dominance results. They will lead to estimates of the transient behavior for rather complex networks, in terms of the known transient behavior for the $M/M/1$ queue.

In Section 4, we present the Jackson network. This is done purely through the operator formalism, augmented by the use of tensor products. For this queueing system and the rest, we modify the usual Kendall notation to denote an $N$-node Jackson network with single servers as $(M/M/1)^N$. In Section 5, we present the class-dependent version of the $M/M/1$ queue, which we shall denote

as $M^C/M^C/1//GD$. In describing this model, we make use of algebraic tools such as free non-abelian semigroups and tensor algebras which are also known as Fock spaces. Finally, in Section 6, all of these mathematical tools come together to construct a very general class-dependent network model that we denote by $(M^C/M^C/1//GD)^N$.

The techniques employed for proving stochastic dominance carry over completely to time-dependent arrival and service rates. In a future paper we shall combine the results for the $M(t)/M(t)/1$ queue in Massey [4] with the stochastic bounds derived in this paper. We can then make statements about non-stationary queueing networks.

## 2. Preliminaries

For any countable set $E$, consider its counting measure. This is a $\sigma$-finite measure on all subsets of $E$. We then set the measure of each subset of $E$ to be its cardinality. For example, each element of $E$ has measure 1. By a trivial use of the Radon–Nikodym theorem, every other $\sigma$-finite signed measure on the power set of $E$ can be uniquely represented by an function integrated against the counting measure. In particular, if we define an $l_1$-norm with respect to the counting measure on these functions, then $l_1(E)$, the functions on $E$ having a finite $l_1$-norm, represent every possible bounded $\sigma$-finite signed measure. The cone of positive functions in $l_1(E)$ with $l_1$-norm equal to unity, represent all of the $\sigma$-finite probability measures on $E$. Given this equivalence, the Kolmogorov forwards equations governing a process of measures on $E$ that evolve in time can be written as

$$\frac{d}{dt}\boldsymbol{p}(t) = \boldsymbol{p}(t)\mathbf{A}$$

where $\boldsymbol{p}(t)$ belongs to the Banach space $l_1(E)$ and $\mathbf{A}$ is a linear operator that acts on $l_1(E)$. We shall refer to $\mathbf{A}$ as the *generator* of the Markov process. For most queueing models, it will be a bounded operator.

We now distinguish some special elements of $l_1(E)$. For any element $\sigma$ of $E$, let $\boldsymbol{e}_\sigma$ be the indicator function for the singleton set $\{\sigma\}$ where

$$\boldsymbol{e}_\sigma(\tau) = \begin{cases} 1, & \tau = \sigma \\ 0, & \tau \neq \sigma. \end{cases}$$

It is clear that $\boldsymbol{e}_\sigma$ belongs to $l_1(E)$. Moreover, $\{\boldsymbol{e}_\sigma\}_{\sigma \in E}$ is a natural basis for $l_1(E)$ since

$$\boldsymbol{f} = \sum_{\sigma \in E} f(\sigma)\boldsymbol{e}_\sigma$$

and $f \equiv 0$ if and only if $f(\sigma) = 0$ for all $\sigma$ in $E$. The cone generated by the $e_\sigma$ coincides with the set of positive vectors in $l_1(E)$. Any operator $\boldsymbol{A}$ on $l_1(E)$ is said to be *positive* if it maps the positive cone into itself. From this notion of positivity, we define two important partial ordering relations. For any $f$ and $g$ in $l_1(E)$, we say that $f \geqq g$ whenever $f - g$ is positive. Similarly, for any two operators $\boldsymbol{A}$ and $\boldsymbol{B}$ on $l_1(E)$, say that $\boldsymbol{A} \geqq \boldsymbol{B}$ if $\boldsymbol{A} - \boldsymbol{B}$ is positive.

For the sake of completeness, we also define $l_\infty(E)$, the set of bounded functions on $E$. Any element $g$ belonging to $l_\infty(E)$ can be written as

$$g = \sum_{\sigma \in E} g(\sigma) e_\sigma$$

where $|g|_\infty = \sup_{\sigma \in E} |g(\sigma)| < \infty$. Of course $l_\infty(E)$ is dual to $l_1(E)$, so any $g$ in $l_\infty(E)$ acts on any $f$ in $l_1(E)$ as

$$f \cdot g = \sum_{\sigma \in E} f(\sigma) g(\sigma).$$

For any bounded operator $\boldsymbol{A}$ on $l_1(E)$, we have $f \boldsymbol{A}$ belonging to $l_1(E)$ and $\boldsymbol{A} g$ belonging to $l_\infty(E)$.

We define $\mathbf{1}$ to be a special $l_\infty(E)$-vector where

$$\mathbf{1} = \sum_{\sigma \in E} e_\sigma.$$

A crucial property for positive $l_1(E)$-vectors is that for $f \geqq 0$, we have

$$|f|_1 = f \cdot \mathbf{1}.$$

Moreover, for positive bounded operators $\boldsymbol{A}$, if $|\boldsymbol{A}|_1$ is the norm induced on $\boldsymbol{A}$ through $l_1(E)$, then

$$|\boldsymbol{A}|_1 = |\boldsymbol{A}\mathbf{1}|_\infty.$$

We now characterize all bounded operators that serve as generators for Markov processes.

*Proposition* 2.1. *For a bounded operator $\boldsymbol{A}$ on $l_1(E)$, the following statements are equivalent:*
(1) $\boldsymbol{A}$ *is the generator for a Markov process on $E$.*
(2) $\boldsymbol{A}\mathbf{1} = 0$ *and* $\boldsymbol{A} + |\boldsymbol{A}|_1 \boldsymbol{I} \geqq 0$.

*Proof.*
(1) $\Rightarrow$ (2). For all probability distributions $\boldsymbol{p}(0)$, we have that $\boldsymbol{p}(0) \exp(t\boldsymbol{A})$ is a probability distribution too. From this it follows that $\boldsymbol{p}(0) \exp(t\boldsymbol{A})\mathbf{1} = \boldsymbol{p}(0) \cdot \mathbf{1} = 1$, hence $\exp(t\boldsymbol{A})\mathbf{1} = \mathbf{1}$. If we differentiate with respect to $t$ and set $t$ equal to 0, then we get $\boldsymbol{A}\mathbf{1} = 0$.

Let $\sigma$ and $\tau$ be two distinct elements in $E$, then $e_\sigma \cdot e_\tau = 0$. If we set $\boldsymbol{p}(0) = e_\sigma$,

then $e_\sigma \exp(tA) \geqq 0$ and so

$$e_\sigma \frac{\exp(tA) - I}{t} e_\tau = e_\sigma \frac{\exp(tA)}{t} e_\tau \geqq 0$$

and as $t$ goes to 0, we get $e_\sigma A e_\tau \geqq 0$. Any diagonal term has the form $e_\sigma A e_\sigma$. Since $|e_\sigma A e_\sigma| \leqq |A|_1$, we have $A + |A|_1 I \geqq 0$.

$(2) \Rightarrow (1)$. We want to show that for $p(0)$ ranging over all probability distributions, we have $|p(0) \exp(tA)|_1 = 1$ and $p(0) \exp(tA) \geqq 0$. Since $\exp(tA) = \exp(t(A + |A|_1 I)) \cdot e^{-|A|_1 t} \geqq 0$, then $p(0) \exp(tA)$ is positive, hence

$$|p(0) \exp(tA)|_1 = p(0) \exp(tA) \cdot 1$$
$$= p(0) \cdot 1$$
$$= 1.$$

The last step follows from the power-series expansion of $\exp(tA)$ and the fact that $A1 = 0$.

## 3. The *M/M/1* queue

The usual formulation of this system goes as follows. Let $Q(t)$ be the queue-length process. The state space $E$ is the set of non-negative integers. The arrivals form a Poisson process with rate $\lambda$ and each customer receives service for a duration exponentially distributed with rate $\mu$. If $p_n(t) = \Pr\{Q(t) = n\}$, then the $p_n(t)$'s solve

$$(3.1) \qquad \frac{d}{dt} p_0(t) = \mu p_1(t) - \lambda p_0(t)$$

and for $n \geqq 1$

$$(3.2) \qquad \frac{d}{dt} p_n(t) = \lambda p_{n-1}(t) + \mu p_{n+1}(t) - (\lambda + \mu) p_n(t).$$

Heuristically, we can interpret these equations as flows in and out of states. For example, in (3.1) and (3.2), $(d/dt)p_n(t)$ is to be the rate of 'flow' through the state $\{Q(t) = n\}$ representing $n$ customers in the system. It then follows that since $\lambda p_{n-1}(t)$ is a positive quantity, it is a flow *into* state $\{Q(t) = n\}$ from state $\{Q(t) = n - 1\}$ via an arrival that occurs at rate $\lambda$. A similar interpretation exists for $\mu p_{n+1}(t)$. The quantity $-(\lambda + \mu)p_n(t)$ is negative, so this represents the rate of flow *out* of the state $\{Q(t) = n\}$. This can be due to an arrival at rate $\lambda$, or a service at rate $\mu$. Since there must always be a non-negative number of customers, the subsequent modifications on the flows are made in (3.1). While these birth and death equations are manageable for the *M/M/1* system, they can become quite tedious for more complex systems, especially when one deals

with the 'boundary behavior' that makes (3.1) differ in form from (3.2). We now rewrite the equations from the point of view of the previous section. Recall that the state space $E$ is the set of non-negative integers. Since the underlying structure of $E$ is totally ordered, we shall simply think of $l_1(E)$ as being a space of sequences and refer to it as $l_1$. Since $\boldsymbol{p}(t)$ belongs to $l_1$, we can write it either as

$$\boldsymbol{p}(t) = \sum_{n=0}^{\infty} p_n(t)\boldsymbol{e}_n$$

or as $\boldsymbol{p}(t) = [p_0(t), p_1(t), p_2(t), \ldots]$. Let $\boldsymbol{R}$ be the right-shift operator on $l_1$ where $\boldsymbol{e}_n\boldsymbol{R} = \boldsymbol{e}_{n+1}$ for all $n \geq 0$. Let $\boldsymbol{L}$ be the left-shift operator where $\boldsymbol{e}_n\boldsymbol{L} = \boldsymbol{e}_{n-1}$ for $n \geq 1$ and $\boldsymbol{e}_0\boldsymbol{L} = 0$. These operators can be seen as fundamental, primitive operators that correspond to the arrival and departure of customers. In fact, we can write $\boldsymbol{A}$ for $\boldsymbol{M/M/1}$ as

(3.3) $$\boldsymbol{A} = \lambda\boldsymbol{R} + \mu\boldsymbol{L} - \lambda\boldsymbol{I} - \mu\boldsymbol{LR}$$

where $\boldsymbol{I}$ is the identity operator. There is a one-to-one correspondence between the various states $\{Q(t) = n\}$ and the unit vectors $\boldsymbol{e}_n$. Given the properties of $\boldsymbol{R}$ and $\boldsymbol{L}$, we see that (3.3) encodes all of the birth and death equations.

For the $\boldsymbol{M/M/1}$ queue, we shall present the two main themes that we shall carry over to more complex systems. First, the verification of the steady-state distribution. We know that the equilibrium distribution for the $\boldsymbol{M/M/1}$ queue (when it exists) is the geometric distribution or

$$\lim_{t\to\infty} \Pr\{Q(t) = n\} = (1-\rho)\rho^n$$

where $\rho < 1$. Now suppose that we encode this distribution as an $l_1$-vector. Let $\boldsymbol{g} = (1-\rho) \cdot \sum_{n=0}^{\infty} \rho^n \boldsymbol{e}_n$ be this representation. We can also write it out as

$$\boldsymbol{g} = (1-\rho)[1, \rho, \rho^2, \cdots].$$

The key algebraic feature $\boldsymbol{g}$ has, is that it is an eigenvector of $\boldsymbol{L}$, the left-shift operator, and $\boldsymbol{gL} = \rho\boldsymbol{g}$. If we substitute this result into $\boldsymbol{gA}$, we get

$$\boldsymbol{gA} = \boldsymbol{g}[\lambda\boldsymbol{R} + \mu\boldsymbol{L} - \lambda\boldsymbol{I} - \mu\boldsymbol{LR}] = \boldsymbol{g}[\lambda\boldsymbol{R} + \mu\rho\boldsymbol{I} - \lambda\boldsymbol{I} - \mu\rho\boldsymbol{R}].$$

We purposely left $\rho$ unspecified, so we may derive from setting $\boldsymbol{gA} = 0$ that $\rho$ must equal $\lambda/\mu$.

The second theme we shall illustrate here is that of stochastic ordering. Consider the set of all non-negative integer-valued random variables. We define a partial ordering on them as follows. For any such $X$ and $Y$, we say that $X \leq_{\text{st}} Y$ if

$$\Pr\{X \geq n\} \leq \Pr\{Y \geq n\}$$

for all positive integers $n$. $X$ is said to be *stochastically dominated* by $Y$. For two **M/M/1** queue-length processes $Q_1(t)$ and $Q_2(t)$, we can give a simple criterion for stochastic dominance in terms of the arrival and service rates $\lambda_1$, $\mu_1$, $\lambda_2$, and $\mu_2$. A special case of a result stated in Kirstein, Franken and Stoyan [3] says that $Q_1(t) \leqq_{st} Q_2(t)$ for all $t \geqq 0$ if $Q_1(0) \leqq_{st} Q_2(0)$, $\lambda_1 \leqq \lambda_2$, and $\mu_1 \geqq \mu_2$. We shall prove this result by operator techniques.

First, we define an operator $\boldsymbol{K}$, where for all $\boldsymbol{e}_n$ we have

$$(3.4) \qquad\qquad \boldsymbol{e}_n \boldsymbol{K} = \sum_{m=0}^{n} \boldsymbol{e}_m.$$

$\boldsymbol{K}$ is a positive operator and we can say formally that

$$\boldsymbol{K} = \boldsymbol{I} + \boldsymbol{L} + \boldsymbol{L}^2 + \cdots = (\boldsymbol{I} - \boldsymbol{L})^{-1}.$$

If $\boldsymbol{p}$, belonging to $l_1$, represents the distribution for some random variable $X$, it is clear by (3.4) that

$$\boldsymbol{p}\boldsymbol{K} = \sum_{n=0}^{\infty} \Pr\{X \geqq n\}\boldsymbol{e}_n.$$

It should be noted that if $E(X) < \infty$, then $|\boldsymbol{p}\boldsymbol{K}|_1 = E(X) + 1$. So all distribution vectors with a finite mean belong to the domain of $\boldsymbol{K}$. Finally, if $\boldsymbol{q}$ represents the distribution for a random variable $Y$, then

$$X \leqq_{st} Y \quad \text{iff} \quad \boldsymbol{p}\boldsymbol{K} \leqq \boldsymbol{q}\boldsymbol{K}.$$

*Lemma* 3.1. *Let* $Q_1(t)$ *and* $Q_2(t)$ *be* **M/M/1** *queue-length processes with the same arrival and service rates,* $\lambda$ *and* $\mu$. *If* $Q_1(0) \leqq_{st} Q_2(0)$, *then* $Q_1(t) \leqq_{st} Q_2(t)$ *for all* $t \geqq 0$.

*Proof.* Let $\boldsymbol{p}$ represent the distribution of $Q_1(0)$, and similarly $\boldsymbol{q}$ for $Q_2(0)$. At time $t$, the distributions for $Q_1(t)$ and $Q_2(t)$ are respectively $\boldsymbol{p} \cdot \exp(\boldsymbol{A}t)$ and $\boldsymbol{q} \cdot \exp(\boldsymbol{A}t)$, where $\boldsymbol{A} = \lambda \boldsymbol{R} + \mu \boldsymbol{L} - \lambda \boldsymbol{I} - \mu \boldsymbol{L}\boldsymbol{R}$. We wish to show that $\boldsymbol{p}\boldsymbol{K} \leqq \boldsymbol{q}\boldsymbol{K}$ implies that $\boldsymbol{p} \exp(t\boldsymbol{A})\boldsymbol{K} \leqq \boldsymbol{q} \exp(t\boldsymbol{A})\boldsymbol{K}$. Since $\boldsymbol{K}^{-1} = (\boldsymbol{I} - \boldsymbol{L})$, it is sufficient to show that $\boldsymbol{K}^{-1} \exp(t\boldsymbol{A})\boldsymbol{K} = \exp(t\boldsymbol{K}^{-1}\boldsymbol{A}\boldsymbol{K}) \geqq 0$. Moreover, by arguments in Section 2, we need only show that there is a positive real number $\alpha$, such that $\boldsymbol{K}^{-1}\boldsymbol{A}\boldsymbol{K} + \alpha \boldsymbol{I} \geqq 0$.

Using the 'power series' representation of $\boldsymbol{K}$, we can readily show that $\boldsymbol{R}\boldsymbol{K} = \boldsymbol{K} + \boldsymbol{R}$ and $\boldsymbol{L}\boldsymbol{K} = \boldsymbol{K} - \boldsymbol{I}$. Therefore,

$$\begin{aligned}
\boldsymbol{A}\boldsymbol{K} &= [\lambda \boldsymbol{R} + \mu \boldsymbol{L} - \lambda \boldsymbol{I} - \mu \boldsymbol{L}\boldsymbol{R}]\boldsymbol{K} \\
&= \lambda(\boldsymbol{K} + \boldsymbol{R}) + \mu \boldsymbol{L}\boldsymbol{K} - \lambda \boldsymbol{K} - \mu(\boldsymbol{L}\boldsymbol{K} + \boldsymbol{L}\boldsymbol{R}) \\
&= \lambda \boldsymbol{R} - \mu \boldsymbol{L}\boldsymbol{R}
\end{aligned}$$

and finally

$$\mathbf{K}^{-1}\mathbf{A}\mathbf{K} = (\mathbf{I} - \mathbf{L})(\lambda \mathbf{R} - \mu \mathbf{L}\mathbf{R})$$
$$= \lambda \mathbf{R} + \mu \mathbf{L}^2 \mathbf{R} - (\lambda + \mu)\mathbf{L}\mathbf{R}$$
$$\geqq -(\lambda + \mu)\mathbf{I}$$

so $\mathbf{K}^{-1}\mathbf{A}\mathbf{K} + (\lambda + \mu)\mathbf{I} \geqq 0$ and this completes the proof.

A Markov process on the positive integers with the above property is said to be *monotone*. This is a crucial property for comparing different processes.

**Proposition 3.2.** *Let* $Q_i(t)$ $(i = 1, 2)$ *be two* **M/M/1** *queues with arrival and service rates* $\lambda_i$ *and* $\mu_i$ *respectively* $(i = 1, 2)$. *If* $\lambda_1 \leqq \lambda_2$, $\mu_1 \geqq \mu_2$, *and* $Q_1(0) \leqq_{st} Q_2(0)$, *then* $Q_1(t) \leqq_{st} Q_2(t)$.

**Proof.** Let $\mathbf{A}_i = \lambda_i \mathbf{R} + \mu_i \mathbf{L} - \lambda_i \mathbf{I} - \mu_i \mathbf{L}\mathbf{R}$ and take $\mathbf{p}_i$ to be the vector representation of the distribution for $Q_i(0)$. We want to show that

$$\mathbf{p}_1 \exp(t\mathbf{A}_1)\mathbf{K} \leqq \mathbf{p}_2 \exp(t\mathbf{A}_2)\mathbf{K}.$$

By the monotonicity of $Q_1(t)$, $\mathbf{p}_1 \exp(t\mathbf{A}_1)\mathbf{K} \leqq \mathbf{p}_2 \exp(t\mathbf{A}_1)\mathbf{K}$, so it is sufficient to show that

$$\exp(t\mathbf{A}_1)\mathbf{K} \leqq \exp(t\mathbf{A}_2)\mathbf{K}.$$

If we apply '$d/ds$' to $\exp(s\mathbf{A}_2)\exp((t-s)\mathbf{A}_1)$, we then get $\exp(s\mathbf{A}_2)$ $(\mathbf{A}_2 - \mathbf{A}_1)\exp((t-s)\mathbf{A}_1)$. Integrating gives us

$(\exp(t\mathbf{A}_2) - \exp(t\mathbf{A}_1))\mathbf{K}$

$$= \int_0^t \exp(s\mathbf{A}_2)(\mathbf{A}_2 - \mathbf{A}_1)\exp((t-s)\mathbf{A}_1)\mathbf{K}\, ds$$

$$= \int_0^t \exp(s\mathbf{A}_2)(\mathbf{A}_2 - \mathbf{A}_1)\mathbf{K} \cdot \mathbf{K}^{-1}\exp((t-s)\mathbf{A}_1)\mathbf{K}\, ds.$$

Now, $(\mathbf{A}_2 - \mathbf{A}_1)\mathbf{K} = (\lambda_2 - \lambda_1)\mathbf{R} + (\mu_1 - \mu_2)\mathbf{L}\mathbf{R}$ is a positive operator by hypothesis. $\mathbf{K}^{-1}\exp((t-s)\mathbf{A}_1)\mathbf{K}$ is positive by the monotonicity of $Q_1(t)$. Finally, $\exp(s\mathbf{A}_2)$ is positive, being the semigroup for a Markov process. We then have $\exp(t\mathbf{A}_2)\mathbf{K} \geqq \exp(t\mathbf{A}_1)\mathbf{K}$.

For future reference, we note the above proof required that $\mathbf{A}_2\mathbf{K} \geqq \mathbf{A}_1\mathbf{K}$ and *only* that $Q_1(t)$ (or $Q_2(t)$) be a monotone process.

The author in Massey [4] used this operator-theoretic approach to do an asymptotic analysis of $\mathbf{M(t)/M(t)/1}$, the time-dependent **M/M/1** queue.

## 4. The $(M/M/1)^N$ network

We now wish to construct a network of $M/M/1$ queues as follows. Given $N$ nodes, let the $i$th node be an $M/M/1$ queue with Poisson arrival rate $\lambda_i$ and exponential service rate $\mu_i$. We 'hook up' this network by an $N \times N$ switching matrix $\boldsymbol{P}$ that is substochastic. A customer having finished service at the $i$th node, arrives at the $j$th node with probability $p_{ij}$. With probability $q_i = 1 - \sum_{j=1}^N p_{ij}$, the customer may decide to leave the total system altogether. We shall assume that $p_{ii} = 0$ for all $i$. This is the Jackson network which we shall denote as $(M/M/1)^N$.

If $Q_i(t)$ is the queue-length process for the $i$th node, then $(Q_1(t), \cdots, Q_N(t))$ is a Markov process. Its state space $E$ equals the set of $N$-tuples $(n_1, \cdots, n_N)$ where each $n_i$ ranges over the non-negative integers. Since we have a natural way of decomposing $E$ into the $N$-fold Cartesian product of the non-negative integers, we want to do a similar decomposition for $l_1(E)$. This can be accomplished through the use of tensor products.

Given two sets $E$ and $D$, let $E \times D$ be their Cartesian product. The set $\{e_{(\sigma,\tau)} \mid \sigma \in E, \tau \in D\}$ is a basis for $l_1(E \times D)$. Recall however that $e_{(\sigma,\tau)}$ is the indicator function for the singleton set $\{(\sigma, \tau)\}$: it then follows that

$$e_{(\sigma,\tau)}(\tilde{\sigma}, \tilde{\tau}) = e_\sigma(\tilde{\sigma}) \cdot e_\tau(\tilde{\tau}).$$

The latter expression has the additional algebraic structure of being linear in $e_\sigma$ and $e_\tau$. We will denote it as $[e_\sigma \otimes e_\tau](\tilde{\sigma}, \tilde{\tau})$. So in a natural way, we can make $l_1(E \times D)$ isomorphic to $l_1(E) \otimes l_1(D)$. For any $\boldsymbol{f}$ in $l_1(E)$ and $\boldsymbol{g}$ in $l_1(D)$, define $\boldsymbol{f} \otimes \boldsymbol{g}$ to be in $l_1(E) \otimes l_1(D)$ and equal to

$$\boldsymbol{f} \otimes \boldsymbol{g} = \sum_{\sigma \in E} \sum_{\tau \in D} f(\sigma)g(\tau)e_\sigma \otimes e_\tau.$$

Furthermore, let the $l_1$-norm on $l_1(E) \otimes l_1(D)$ be the same as that on $l_1(E \times D)$. Consequently $|\boldsymbol{f} \otimes \boldsymbol{g}|_1 = |\boldsymbol{f}|_1 |\boldsymbol{g}|_1$, for all given $\boldsymbol{f}$ and $\boldsymbol{g}$. We list the bilinearity properties that $\boldsymbol{f} \otimes \boldsymbol{g}$ has:

$$(\boldsymbol{f}_1 + \boldsymbol{f}_2) \otimes \boldsymbol{g} = \boldsymbol{f}_1 \otimes \boldsymbol{g} + \boldsymbol{f}_2 \otimes \boldsymbol{g}$$
$$\boldsymbol{f} \otimes (\boldsymbol{g}_1 + \boldsymbol{g}_2) = \boldsymbol{f} \otimes \boldsymbol{g}_1 + \boldsymbol{f} \otimes \boldsymbol{g}_2$$
$$(\alpha \boldsymbol{f}) \otimes \boldsymbol{g} = \boldsymbol{f} \otimes (\alpha \boldsymbol{g}) = \alpha(\boldsymbol{f} \otimes \boldsymbol{g}).$$

The $\{e_\sigma \otimes e_\tau\}_{\sigma \in E, \tau \in D}$ are a basis for $l_1(E) \otimes l_1(D)$, so to define an operator on this space it is sufficient to define a linear operator on these basis elements. For example, if $\boldsymbol{A}$ is an operator on $l_1(E)$ and $\boldsymbol{B}$ is one on $l_1(D)$, then $\boldsymbol{A} \otimes \boldsymbol{B}$ is defined on $l_1(E) \otimes l_1(D)$ by setting

$$(e_\sigma \otimes e_\tau)[\boldsymbol{A} \otimes \boldsymbol{B}] = e_\sigma \boldsymbol{A} \otimes e_\tau \boldsymbol{B}$$

for each basis element.

It is clear how the above constructions can be extended to a tensoring of $N$ spaces together where $N$ exceeds 2. We then see that given the state space $E$ for $(M/M/1)^N$ gives us an isomorphism between $l_1(E)$ and $l_1^{(N)}$, which is $l_1$ tensored with itself $N$ times. For the $M/M/1$ queue, $R$ and $L$ were natural operators to define on $l_1$. Now we shall use them to define our basic, primitive operators for $(M/M/1)^N$. We define $R_i$ and $L_i$ for $i = 1, \cdots, N$ such that

$$R_i = I \otimes \cdots \otimes R \otimes \cdots \otimes I \quad (i\text{th place})$$

$$L_i = I \otimes \cdots \otimes L \otimes \cdots \otimes I \quad (i\text{th place}).$$

These operators now enable us to construct the generator $A$ for the $(M/M/1)^N$ network

(4.1) $$A = \sum_{i=1}^{N} \left[ \lambda_i R_i + \mu_i q_i L_i + \sum_{j=1}^{N} \mu_i p_{ij} L_i R_j - \lambda_i I - \mu_i L_i R_i \right].$$

Paralleling the interplay between the primitive operators and the 'flows' in the birth and death equations for the $M/M/1$ queue, we can verify that this is the correct generator. A state $(n_1, \cdots, n_N)$ is represented by the basis vector $e_{n_1} \otimes \cdots \otimes e_{n_N}$. For simplicity, let $i = 1$ and $j = N$. A customer may arrive to node 1 with arrival rate $\lambda_1$, adding a customer to the queue. This corresponds to

$$(e_{n_1} \otimes \cdots \otimes e_{n_N}) R_1 = e_{n_1+1} \otimes \cdots \otimes e_{n_N}.$$

If a customer exists at the first node then service will occur with rate $\mu_1$, and with probability $q_1$ the customer will leave the entire network. We encode this as

$$(e_{n_1} \otimes \cdots \otimes e_{n_N}) L_1 = e_{n_1-1} \otimes \cdots \otimes e_{n_N}$$

if $n_1 \geqq 1$, otherwise we get the zero vector. Finally, this departing customer may instead decide with probability $p_{1N}$ to transfer and enter the $N$th node. This is interpreted as

$$(e_{n_1} \otimes \cdots \otimes e_{n_N}) L_1 R_N = e_{n_1-1} \otimes \cdots \otimes e_{n_N+1},$$

again if $n_1 \geqq 1$. A similar description can be given for the negative terms.

Given this formulation, we can give a purely algebraic proof of Jackson's theorem, assuming only that the steady-state distribution has a product form. If we let $g_i$, belonging to $l_1$, represent the marginal equilibrium distribution for $Q_i(t)$, then a product-form solution means that

$$\left( \bigotimes_{i=1}^{N} g_i \right) A = 0.$$

The rank-$N$ tensor $\bigotimes_{i=1}^{N} g_i$ has the property that summing on all indices except

for the $i$th one, gives back the $i$th vector in the original tensor product, up to a scalar. Moreover, the primitive operators $\boldsymbol{L}_i$ and $\boldsymbol{R}_i$ preserve the tensor-product form when applied to $\otimes_{i=1}^{N} \boldsymbol{g}_i$. Consequently, if we sum on all but the $i$th index of the rank-$N$ tensor $(\otimes_{i=1}^{N} \boldsymbol{g}_i)\boldsymbol{A}$, for some positive scalars $\xi_i$ and $\eta_i$ we have

$$\boldsymbol{g}_i(\xi_i \boldsymbol{R} + \eta_i \boldsymbol{L} - \xi_i \boldsymbol{I} - \eta_i \boldsymbol{L}\boldsymbol{R}) = 0.$$

Hence $\boldsymbol{g}_i$ must be the steady-state distribution for an $\boldsymbol{M/M/1}$ queue, so there is some $0 < \rho_i < 1$ such that $\boldsymbol{g}_i \boldsymbol{L} = \rho_i \boldsymbol{g}_i$ for each $i$. From this, it follows that for each $j$,

$$\left( \overset{N}{\underset{i=1}{\bigotimes}} \boldsymbol{g}_i \right) \boldsymbol{L}_j = \rho_j \overset{N}{\underset{i=1}{\bigotimes}} \boldsymbol{g}_i.$$

Let $\theta_i = \rho_i \mu_i$; substituting into $(\otimes_{i=1}^{N} \boldsymbol{g}_i)\boldsymbol{A}$ gives

$$\left( \overset{N}{\underset{i=1}{\bigotimes}} \boldsymbol{g}_i \right) \boldsymbol{A} = \left( \overset{N}{\underset{i=1}{\bigotimes}} \boldsymbol{g}_i \right) \sum_{i=1}^{N} \left[ \lambda_i \boldsymbol{R}_i + \mu_i \rho_i q_i \boldsymbol{I} + \sum_{j=1}^{N} \mu_i \rho_i p_{ij} \boldsymbol{R}_j - \lambda_i \boldsymbol{I} - \mu_i \rho_i \boldsymbol{R}_i \right]$$

$$= \left( \overset{N}{\underset{i=1}{\bigotimes}} \boldsymbol{g}_i \right) \sum_{i=1}^{N} \left[ \left( \lambda_i + \sum_{j=1}^{N} \theta_j p_{ji} - \theta_i \right) \boldsymbol{R}_i + (\theta_i q_i - \lambda_i) \boldsymbol{I} \right].$$

If we want $(\otimes_{i=1}^{N} \boldsymbol{g}_i)\boldsymbol{A} = 0$, then we must have for each $i$, $\theta_i < \mu_i$

$$\theta_i = \sum_{j=1}^{N} \theta_j p_{ji} + \lambda_i$$

and from this follows $\sum_{i=1}^{N} \lambda_i = \sum_{i=1}^{N} \theta_i q_i$.

Aside from reproving a well-known result in a less *ad hoc* manner, we can put this machinery to work to prove new results. In Massey [5], the author derived bounds for the transient behavior of the $(\boldsymbol{M/M/1})^N$ network in terms of the known transient behavior for the $\boldsymbol{M/M/1}$ queue. For a given $(\boldsymbol{M/M/1})^N$ network, let $(X_1(t), \cdots, X_N(t))$ be a collection of independent $\boldsymbol{M/M/1}$ queue-length processes with $X_i(0) = Q_i(0)$, arrival rate $\lambda_i + \sum_{j=1}^{N} \mu_j p_{ji}$ and service rate $\mu_i$, then for all $t \geq 0$

$$\Pr\{Q_1(t) \geq n_1, \cdots, Q_N(t) \geq n_N\} \leq \prod_{i=1}^{N} \Pr\{X_i(t) \geq n_i\}$$

for all non-negative integers $n_1, \cdots, n_N$. Moreover, if $Y(t)$ is an $\boldsymbol{M/M/1}$ queue-length process with arrival rate $\sum_{i=1}^{N} \lambda_i$ and service rate $\sum_{i=1}^{N} \mu_i q_i$, then for all $t \geq 0$,

$$(4.2) \qquad \Pr\left\{ \sum_{i=1}^{N} Q_i(t) \geq n \right\} \geq \Pr\{Y(t) \geq n\}$$

for all non-negative integers $n$. We shall briefly rederive the latter of these two

results. Doing so allows us to introduce an aggregation operator $S$ that maps $l_1^{(N)}$ into $l_1$ where if $m = \sum_{i=1}^{N} n_i$, we have

$$\left(\bigotimes_{i=1}^{N} e_{n_i}\right) S = e_m.$$

The distribution for $\sum_{i=1}^{N} Q_i(t)$ can then be represented by $\boldsymbol{p}(0) \exp(t\boldsymbol{A})\boldsymbol{S}$. $\boldsymbol{S}$ is a positive operator, and it has special algebraic properties $\boldsymbol{R}_i\boldsymbol{S} = \boldsymbol{S}\boldsymbol{R}$ and $\boldsymbol{L}_i\boldsymbol{S} \leqq \boldsymbol{S}\boldsymbol{L}$. Now let $\boldsymbol{B} = \lambda'\boldsymbol{R} + \mu'\boldsymbol{L} - \mu'\boldsymbol{LR}$ be the generator for $Y(t)$, where $\lambda' = \sum_{i=1}^{N} \lambda_i$ and $\mu' = \sum_{i=1}^{N} \mu_i q_i$. Inequality (4.2) is then equivalent to

$$\boldsymbol{p}(0) \exp(t\boldsymbol{A})\boldsymbol{SK} \geqq \boldsymbol{p}(0)\boldsymbol{S} \exp(t\boldsymbol{B})\boldsymbol{K}.$$

So it is sufficient to show that

$$\boldsymbol{S} \exp(t\boldsymbol{B})\boldsymbol{K} \leqq \exp(t\boldsymbol{A})\boldsymbol{SK}.$$

Using the same argument as in Proposition 3.2, we differentiate $\exp(s\boldsymbol{A})\boldsymbol{S} \exp((t-s)\boldsymbol{B})\boldsymbol{K}$ to deduce that

$\exp(t\boldsymbol{A})\boldsymbol{SK} - \boldsymbol{S} \exp(t\boldsymbol{B})\boldsymbol{K}$

$$= \int_0^t \exp(s\boldsymbol{A})(\boldsymbol{AS} - \boldsymbol{SB}) \exp((t-s)\boldsymbol{B})\boldsymbol{K} \, ds$$

$$= \int_0^t \exp(s\boldsymbol{A})(\boldsymbol{AS} - \boldsymbol{SB})\boldsymbol{K} \cdot \boldsymbol{K}^{-1} \exp((t-s)\boldsymbol{B})\boldsymbol{K} \, ds.$$

$Y(t)$ is monotone, so $\boldsymbol{K}^{-1} \exp((t-s)\boldsymbol{B})\boldsymbol{K} \geqq 0$. So we need only show that $\boldsymbol{ASK} \geqq \boldsymbol{SBK}$. Using the algebraic properties of $\boldsymbol{S}$ and $\boldsymbol{K}$ gives

$$\boldsymbol{ASK} = \sum_{i=1}^{N} \left[ \lambda_i\boldsymbol{R}_i + \mu_i q_i\boldsymbol{L}_i + \sum_{j=1}^{N} \mu_i p_{ij}\boldsymbol{L}_i\boldsymbol{R}_j - \lambda_i\boldsymbol{I} - \mu_i\boldsymbol{L}_i\boldsymbol{R}_i \right]\boldsymbol{SK}$$

$$= \sum_{i=1}^{N} \left[ \lambda_i\boldsymbol{SR} + \mu_i q_i\boldsymbol{L}_i\boldsymbol{S} + \sum_{j=1}^{N} \mu_i p_{ij}\boldsymbol{L}_i\boldsymbol{SR} - \lambda_i\boldsymbol{S} - \mu_i\boldsymbol{L}_i\boldsymbol{SR} \right]\boldsymbol{K}$$

$$= \sum_{i=1}^{N} [\lambda_i\boldsymbol{S}(\boldsymbol{R} - \boldsymbol{I}) + \mu_i q_i\boldsymbol{L}_i\boldsymbol{S}(\boldsymbol{I} - \boldsymbol{R})]\boldsymbol{K}$$

$$= \sum_{i=1}^{N} [\lambda_i\boldsymbol{SR} - \mu_i q_i\boldsymbol{L}_i\boldsymbol{SR}]$$

$$\geqq \boldsymbol{S}[\lambda'\boldsymbol{R} - \mu'\boldsymbol{LR}]$$

$$\geqq \boldsymbol{SBK}$$

and this proves the result.

In Section 6, we shall generalize these results for $(\boldsymbol{M}^C/\boldsymbol{M}^C/1//\boldsymbol{GD})^{\boldsymbol{N}}$ networks. The product-form equilibrium distribution carries over to a special subclass of these networks first derived by Kelly [2]. However, the bounds for

estimating the transient behavior can be derived for any $(M^C/M^C/1//GD)^N$ network.

## 5. The $M^C/M^C/1//GD$ queue

We now wish to construct a very general single-server queue that allows for different classes of customer. Let $C$ be a finite set of 'tags', one for each class. Representative elements of $C$ will be denoted by $\alpha$ or $\beta$. For the $\alpha$-class customers, $\lambda_\alpha$ is their Poisson arrival rate and $\mu_\alpha$ is their exponential service rate. With each new service, one must be aware of the class of the customer to be served. The state space $E$ must list all possible configurations for queueing lines. By this we mean the number of customers in line as well as the class of each customer in line. Moreover, when we insert customers in line or delete them, we do so by a discipline. For example, suppose we want to delete a customer from a configuration $\sigma$ belonging to $E$. Let $|\sigma|$ be the number of customers in line. If the first position is the head of the line, we shall delete a customer in the $i$th position with probability $\psi_i(\sigma)$ where $i = 1, \cdots, |\sigma|$. So for each $\sigma$ in $E$ we have $\sum_{i=1}^{|\sigma|} \psi_i(\sigma) = 1$. It is this family of probability distributions that we shall call a *deletion discipline*. We shall refer to this family as $\Psi$ which stands for the set $\{\psi_i(\sigma) \mid \sigma \in E \text{ and } i = 1, \cdots, |\sigma|\}$.

We shall distinguish three special disciplines:

(i) $\Psi_F$ deletes customers at the head of the line where

$$\psi_i(\sigma) = \begin{cases} 1, & i = 1 \\ 0, & \text{otherwise,} \end{cases}$$

(ii) $\Psi_L$ deletes customers at the end of the line, and

$$\psi_i(\sigma) = \begin{cases} 1, & i = |\sigma| \\ 0, & \text{otherwise.} \end{cases}$$

(iii) $\Psi_P$ deletes any customer in line with uniform probability, hence $\psi_i(\sigma) = 1/|\sigma|$ for $|\sigma| \neq 0$.

Similarly, we can define an *insertion discipline*, which we denote by $\Phi$, in the same fashion as $\Psi$. However, when we insert an $\alpha$-class customer into a configuration $\sigma$, we regard $(\sigma, \alpha)$ as the new configuration and let $i$ in $\phi_i(\sigma, \alpha)$ range from 1 to $|\sigma| + 1$. We say that $\Phi \simeq \Psi$ if $\tau = (\sigma, \alpha)$ and $\phi_i(\sigma, \alpha) = \psi_i(\tau)$ for all $\alpha$, $\sigma$, and $i$. Thus $\Phi_F$, where $\Phi_F \simeq \Psi_F$, is the discipline for inserting customers at the head of line. We then define $\Phi_L$ and $\Phi_P$ in a similar fashion.

In Kendall notation, we let GD refer to a *general discipline*. We shall define this to be the pair $(\Phi, \Psi)$, where $\Phi$ and $\Psi$ are arbitrary. In the sense that GD $= (\Phi, \Psi)$, we also have FIFO $= (\Phi_L, \Psi_F)$, LIFO $= (\Phi_L, \Psi_L)$, and PS $= (\Phi_L, \Psi_P)$

the processor-sharing discipline. Restricting the disciplines to being *load dependent*, that is $\phi_i(\sigma, \alpha) = \phi_i(|\sigma| + 1)$ and $\psi_i(\sigma) = \psi_i(|\sigma|)$, we then have the type of disciplines as defined in Kelly [2]. Notice that the FIFO, LIFO, and PS disciplines are all of this type.

To describe a service more complicated than a once only exponential holding time, we add to this model the notion of *internal switching*. After a customer of class $\alpha$ completes service, the subsequent departure may not be allowed. With probability $p_{\alpha\beta}$, the customer merely remains in his same spot in line and is transformed into a $\beta$ class customer, and ready to resume service as such. The collection of $p_{\alpha\beta}$'s for $\alpha$ and $\beta$ in $C$ is said to be a $|C| \times |C|$ switching matrix. It is a substochastic one, so with probability $q_\alpha = 1 - \sum_{\beta \in C} p_{\alpha\beta}$, a class $\alpha$ customer will be allowed to leave the queue. For simplicity, we assume that $p_{\alpha\alpha} = 0$ for all $\alpha$.

Consider a FIFO discipline with $C = \{\alpha, \beta\}$, two elements. Let $\lambda_\alpha$ be arbitrary, and set $\lambda_\beta = 0$, $\mu_\alpha = \mu_\beta = \mu$, $p_{\alpha\alpha} = p_{\beta\beta} = p_{\beta\alpha} = 0$, with $p_{\alpha\beta} = 1$. The switching matrix tells us that every class $\alpha$ customer after ending service with rate $\mu$, stays in the server to become a class $\beta$ customer with the same service requirement. The class $\beta$ customer upon completion of service, is told to leave. If we treat the $\alpha$, $\beta$ pair as one customer, then this models an $M/E_2/1//FIFO$ system. $E_2$ is the distribution for the sum of two i.i.d. exponential random variables. Now let $\mu_\alpha \neq \mu_\beta$ and $p_{\alpha\alpha} = p_{\beta\beta} = p_{\alpha\beta} = p_{\beta\alpha} = 0$. Since we can superimpose independent Poisson processes, this models an $M/H_2/1//FIFO$ system. The arrival rate is $\lambda_\alpha + \lambda_\beta$ and the service distribution $H_2$, is a convex combination of two independent exponential distributions where with probability $\dfrac{\lambda_\alpha}{\lambda_\alpha + \lambda_\beta} \left(\text{or } \dfrac{\lambda_\beta}{\lambda_\alpha + \lambda_\beta}\right)$, a customer receives service with mean $\dfrac{1}{\mu_\alpha} \left(\text{or } \dfrac{1}{\mu_\beta}, \text{respectively}\right)$. In a similar fashion, we can use internal switching to model any phase-type service distribution.

We now wish to encode this queueing system into operator notation. First, we must impose some algebraic structure on $E$. For the $M/M/1$ queue, $E$ was the non-negative integers. We made transitions by adding or subtracting one customer. We generalize this for $M^C/M^C/1//GD$ by letting $E = S_C$, the free non-abelian semigroup with identity element, generated on the set $C$. If $C$ is a singleton set, then $S_C$ is isomorphic to the non-negative integers and our model reduces to the $M/M/1$ queue. Let $\oplus$ denote the semigroup operation and let $O$ be the identity element. Every element $\sigma$ of $E$ can then be written as $\alpha_1 \oplus \cdots \oplus \alpha_n$. It is clear that $\alpha_1 \oplus \cdots \oplus \alpha_n$ encodes the state of having $n$ customers in line with the $i$th customer being of the $\alpha_i$ class.

Our forwards equation for this queueing model exists then on the Banach space $l_1(S_C)$. To induce an algebraic structure on $l_1(S_C)$, we modify the natural

map from $\sigma$ in $S_C$ to $\boldsymbol{e}_\sigma$ in $l_1(S_C)$. We modify it to be a semigroup homomorphism where

$$\boldsymbol{e}_{\alpha_1 \oplus \cdots \oplus \alpha_n} = \boldsymbol{e}_{\alpha_1} \otimes \cdots \otimes \boldsymbol{e}_{\alpha_n}$$

and $\otimes$ is the tensor product operation. Thus there is a one-to-one correspondence between the basis vectors for strings of length $n$ and the basis vectors for the space of $|C|$-dimensional tensors of rank $n$. Therefore $l_1(S_C)$ is isomorphic to $\mathcal{F}(l_1(C))$, the *Fock space* of $l_1(C)$ which is the direct sum of the spaces of $|C|$-dimensional tensors of all ranks. This includes the scalars, where we associate $\boldsymbol{e}_0$ with the scalar 1. $\mathcal{F}(l_1(C))$ is also referred to as the *tensor algebra* of $l_1(C)$. An element of $\mathcal{F}(l_1(C))$ can then be thought of as an infinite-dimensional vector, where the $n$th component is a $|C|$-dimensional tensor of rank $n$. In this manner, we can still define the notions of right- and left-shift operators on $\mathcal{F}(l_1(C))$.

We define a generalized right-shift operator as follows. Given a class $\alpha$, and an insertion discipline $\Phi$, we have

$$\left(\overset{n}{\underset{i=1}{\otimes}} \boldsymbol{e}_{\alpha_i}\right) \boldsymbol{R}(\alpha, \Phi) = \sum_{i=1}^{n+1} \phi_i(\sigma \oplus \alpha) \left(\underset{j<i}{\otimes} \boldsymbol{e}_{\alpha_j}\right) \otimes \boldsymbol{e}_\alpha \otimes \left(\underset{j \geq i}{\otimes} \boldsymbol{e}_{\alpha_j}\right)$$

where $\sigma = \alpha_1 \oplus \cdots \oplus \alpha_n$. Just as $\boldsymbol{R}$ maps $\boldsymbol{e}_n$ to $\boldsymbol{e}_{n+1}$, $\boldsymbol{R}(\alpha, \Phi)$ maps rank $n$ tensors into rank $n+1$ tensors. Similarly, given a class $\alpha$, and a deletion discipline $\Psi$, we define a generalized left-shift operator $\boldsymbol{L}(\alpha, \Psi)$ where

$$\left(\overset{n}{\underset{i=1}{\otimes}} \boldsymbol{e}_{\alpha_i}\right) \boldsymbol{L}(\alpha, \Psi) = \sum_{i=1}^{n} \psi_i(\sigma) e_\alpha(\alpha_i) \underset{j \neq i}{\otimes} \boldsymbol{e}_{\alpha_j}$$

and $\sigma = \alpha_1 \oplus \cdots \oplus \alpha_n$. Recall that $e_\alpha(\alpha_i)$ is a scalar, equaling 1 if $\alpha = \alpha_i$ and 0 otherwise. Finally, we must introduce a modification operator $\boldsymbol{M}(\alpha, \beta, \Psi)$, that deletes an $\alpha$-class customer (if possible) via $\Psi$ and then substitutes a $\beta$-class customer in the same place. In other words

$$\left(\overset{n}{\underset{i=1}{\otimes}} \boldsymbol{e}_{\alpha_i}\right) \boldsymbol{M}(\alpha, \beta, \Psi) = \sum_{i=1}^{n} \psi_i(\sigma) e_\alpha(\alpha_i) \left(\underset{j<i}{\otimes} \boldsymbol{e}_{\alpha_j}\right) \otimes \boldsymbol{e}_\beta \otimes \left(\underset{j>i}{\otimes} \boldsymbol{e}_{\alpha_j}\right).$$

We can now construct the generator $\boldsymbol{A}$ for an $\boldsymbol{M}^C/\boldsymbol{M}^C/1//\boldsymbol{GD}$ system:

(5.1)

$$\boldsymbol{A} = \sum_{\alpha \in C} \left[ \lambda_\alpha \boldsymbol{R}(\alpha, \Phi) + \mu_\alpha q_\alpha \boldsymbol{L}(\alpha, \Psi) + \sum_{\beta \in C} \mu_\alpha p_{\alpha\beta} \boldsymbol{M}(\alpha, \beta, \Psi) - \lambda_\alpha \boldsymbol{I} - \mu_\alpha \boldsymbol{M}(\alpha, \Psi) \right]$$

where $\boldsymbol{M}(\alpha, \Psi) = \boldsymbol{M}(\alpha, \alpha, \Psi)$.

Let $\Gamma$ belong to $\mathcal{F}(l_1(C))$. We say that $\Gamma$ represents a *tensor geometric distribution* on $S_C$, if there exists a positive $\boldsymbol{g}$ in $l_1(C)$ with $|\boldsymbol{g}|_1 < 1$ such that

$$\Gamma = (1 - |\boldsymbol{g}|_1)[1, \boldsymbol{g}, \boldsymbol{g} \otimes \boldsymbol{g}, \cdots].$$

The tensor geometric distribution is the steady-state distribution for a special class of $M^C/M^C/1//GD$ queues, as we shall soon see.

We say that $\Psi$ (or $\Phi$) is an *abelian discipline* if for all strings of length $n$ and for all permutations $\pi$ on $n$ objects we have

$$\psi_i(\alpha_1 \oplus \cdots \oplus \alpha_n) = \psi_i(\alpha_{\pi(1)} \oplus \cdots \oplus \alpha_{\pi(n)}).$$

Notice that load-dependent disciplines are a special case of abelian disciplines. Disciplines of this type have more than just a theoretical interest. We shall illustrate here an example of a non-trivial abelian discipline as discussed in Fayolle, Iasnogorodski and Mitrani [1].

For each class $\alpha$, if we assign a positive weight $g(\alpha) > 0$, then we can define the following service discipline:

$$\psi_i(\alpha_1 \oplus \cdots \oplus \alpha_n) = \frac{g(\alpha_i)}{\sum\limits_{\alpha \in C} g(\alpha) |\sigma|_\alpha}$$

where $|\sigma|_\alpha$ is the number of occurrences of $\alpha$ in $\sigma = \alpha_1 \oplus \cdots \oplus \alpha_n$. Notice that if all of the $g(\alpha)$'s are the same, then $\Psi$ reduces to $\Psi_p$. $\Psi$ can then be thought of as a weighted processor-sharing service discipline.

We now give an example of a non-abelian service discipline. Suppose that $C$ is a totally ordered finite set. In other words, for all $\alpha$ and $\beta$ in $C$, we have $\alpha \leq \beta$ or $\beta \leq \alpha$. Now define $\Psi$ as

$$\psi_i(\alpha_1 \oplus \cdots \oplus \alpha_n) = \begin{cases} 1, & i = \inf\{j \mid \alpha_j \geq \alpha_k \text{ for } k = 1, \cdots, n\} \\ 0, & \text{otherwise.} \end{cases}$$

This models a priority service discipline, where service is given to the first in line of the highest-priority class in the queue. With minor modifications, we can model preemptive and non-preemptive priority systems equally well. This illustrates the main purpose of defining disciplines in such generality. Results that hold for any $M^C/M^C/1//GD$ system, then hold for any discipline that is FIFO, LIFO, PS, or even priority scheduling.

*Proposition 5.1. Let $\Gamma$ be a tensor geometric distribution, and let $\Psi$ be an abelian discipline, then*
  (1) $\Gamma \cdot L(\alpha, \Psi) = g(\alpha)\Gamma$,
  (2) $\Gamma \cdot M(\alpha, \Psi) = g(\alpha)\Gamma \cdot R(\alpha, \Phi)$ *where* $\Phi \simeq \Psi$.
*Moreover, if $\Psi$ is merely a load-dependent discipline, then*
  (3) $\Gamma \cdot M(\alpha, \beta, \Psi) = g(\alpha)\Gamma \cdot R(\beta, \Phi)$ *where* $\Phi \simeq \Psi$,
  (4) $\sum_{\alpha \in C} \Gamma \cdot M(\alpha, \Psi)$ *is independent of* $\Psi$.

*Proof.* Given $g$, a $|C|$-dimensional vector, let $g^{(n)}$ be a rank-$n$ tensor equal to $g$ tensored with itself $n$ times. To prove (1), it is sufficient to show that

$g^{(n+1)}\boldsymbol{L}(\alpha, \Psi) = g(\alpha)g^{(n)}$ hence

$$g^{(n+1)}\boldsymbol{L}(\alpha, \Psi) = \sum_{(\alpha_1,\cdots,\alpha_{n+1})\in C^{n+1}} \left(\prod_{i=1}^{n+1} g(\alpha_i)\right)\left(\bigotimes_{i=1}^{n+1} \boldsymbol{e}_{\alpha_i}\right)\boldsymbol{L}(\alpha, \Psi)$$

$$= \sum_{(\alpha_1,\cdots,\alpha_{n+1})\in C^{n+1}} \left(\prod_{i=1}^{n+1} g(\alpha_i)\right) \sum_{j=1}^{n+1} \psi_j(\alpha_1, \cdots, \alpha_{n+1})\boldsymbol{e}_\alpha(\alpha_j) \bigotimes_{k\neq j} \boldsymbol{e}_{\alpha_k}$$

$$= \sum_{j=1}^{n+1} \sum_{(\alpha_1,\cdots,\alpha_{n+1})\in C^{n+1}} \left(\prod_{i=1}^{n+1} g(\alpha_i)\right)\psi_j(\alpha_1, \cdots, \alpha_{n+1})\boldsymbol{e}_\alpha(\alpha_j) \bigotimes_{k\neq j} \boldsymbol{e}_{\alpha_k}$$

$$= \sum_{j=1}^{n+1} g(\alpha) \sum_{(\cdots,\alpha_j,\cdots)\in C^n} \left(\prod_{i\neq j} g(\alpha_j)\right)\psi_j(\alpha_1, \cdots, \alpha_{n+1}, \alpha) \bigotimes_{k\neq j} \boldsymbol{e}_{\alpha_k}$$

$$= g(\alpha) \sum_{(\alpha_1,\cdots,\alpha_n)\in C^n} \left(\prod_{i=1}^{n} g(\alpha_i)\right)\left(\sum_{j=1}^{n+1} \psi_j(\alpha_1, \cdots, \alpha_{n+1}, \alpha)\right) \bigotimes_{k=1}^{n} \boldsymbol{e}_{\alpha_k}$$

$$= g(\alpha) \cdot \boldsymbol{g}^{(n)}.$$

We henceforth abuse notation by setting $\boldsymbol{R}(\alpha, \Psi)$ equal to $\boldsymbol{R}(\alpha, \Phi)$ provided $\Phi \approx \Psi$, and similarly defining $\boldsymbol{M}(\alpha, \beta, \Phi)$ or $\boldsymbol{L}(\alpha, \Phi)$. To prove (2), we need only show the result below.

$$g^{(n+1)}\boldsymbol{M}(\alpha, \Psi)$$

$$= \sum_{(\alpha_1,\cdots,\alpha_{n+1})\in C^{n+1}} \left(\prod_{i=1}^{n+1} g(\alpha_i)\right) \cdot \sum_{j=1}^{n+1} \psi_j(\alpha_1, \cdots, \alpha_{n+1})\boldsymbol{e}_\alpha(\alpha_j)\left(\bigotimes_{k<j} \boldsymbol{e}_{\alpha_k}\right)\otimes \boldsymbol{e}_\alpha \otimes \left(\bigotimes_{k>j} \boldsymbol{e}_{\alpha_k}\right)$$

$$= g(\alpha) \sum_{(\alpha_1,\cdots,\alpha_n)\in C^n} \left(\prod_{i=1}^{n} g(\alpha_i)\right) \sum_{j=1}^{n+1} \psi_j(\alpha_1, \cdots, \alpha_n, \alpha)\left(\bigotimes_{k<j} \boldsymbol{e}_{\alpha_k}\right)\otimes \boldsymbol{e}_\alpha \otimes \left(\bigotimes_{k\geq j} \boldsymbol{e}_{\alpha_k}\right)$$

$$= g(\alpha)\left[\sum_{(\alpha_1,\cdots,\alpha_n)\in C^n} \left(\prod_{i=1}^{n} g(\alpha_j)\right) \bigotimes_{k=1}^{n} \boldsymbol{e}_{\alpha_k}\right]\boldsymbol{R}(\alpha, \Psi)$$

$$= g(\alpha)\boldsymbol{g}^{(n)}\boldsymbol{R}(\alpha, \Psi).$$

When $\Psi$ is merely load dependent, then the $\psi_j$'s have no class dependence. We can then use the above argument to show that $\boldsymbol{g}^{(n+1)}\boldsymbol{M}(\alpha, \beta, \Psi) = g(\alpha)\boldsymbol{g}^{(n)}\boldsymbol{R}(\beta, \Psi)$. Moreover,

$$\sum_{\alpha\in C} \boldsymbol{g}^{(n+1)}\boldsymbol{M}(\alpha, \Psi)$$

$$= \sum_{(\alpha_1,\cdots,\alpha_{n+1})\in C^{n+1}} \left(\prod_{i=1}^{n+1} g(\alpha_i)\right) \sum_{j=1}^{n+1} \psi_j(n+1) \bigotimes_{k=1}^{n+1} \boldsymbol{e}_{\alpha_k}$$

$$= \boldsymbol{g}^{(n+1)}$$

and so we are done.

*Theorem 5.2. For an $\boldsymbol{M}^C/\boldsymbol{M}^C/1//\boldsymbol{GD}$ queue, take $\boldsymbol{A}$ to be its generator and let*

$\Gamma$ *represent a tensor geometric distribution generated by* $\boldsymbol{g} = \sum_{\alpha \in C} \dfrac{\theta_\alpha}{\mu_\alpha} \boldsymbol{e}_\alpha$, *where the* $\theta_\alpha$*'s are to be determined. Suppose that one of the following conditions holds:*

(1) $\Phi$ *and* $\Psi$ *are abelian disciplines with* $\Phi \simeq \Psi$ *and* $p_{\alpha,\beta} = 0$ *for all* $\alpha$ *and* $\beta$ *in C.*

(2) $\Phi$ *and* $\Psi$ *are load-dependent disciplines with* $\Phi \simeq \Psi$.

*If a steady-state distribution for* $\boldsymbol{A}$ *exists, then* $\Gamma \boldsymbol{A} = 0$ *where the* $\theta(\alpha)$*'s solve*

$$(5.2) \qquad \theta_\alpha = \lambda_\alpha + \sum_{\beta \in C} \theta_\beta p_{\beta\alpha}$$

*and*

$$\sum_{\alpha \in C} \frac{\theta_\alpha}{\mu_\alpha} < 1.$$

*Furthermore, suppose instead of conditions* (1) *or* (2), *we have:*

(3) $\Phi$ *and* $\Psi$ *are load-dependent with* $\Phi \neq \Psi$.

*Then* $\Gamma \boldsymbol{A} = 0$ *when we constrain the* $\mu_\alpha$*'s such that there exists some constant* $\kappa$ *where for all* $\alpha$, $\mu_\alpha = \dfrac{\theta_\alpha}{\lambda_\alpha} \cdot \kappa$ *and the* $\theta_\alpha$*'s solve* (5.2).

*Proof.* Using the results of Proposition 5.1 and hypothesis (1) or (2) gives us

$$\Gamma \boldsymbol{A} = \Gamma \cdot \sum_{\alpha \in C} \left[ \lambda_\alpha \boldsymbol{R}(\alpha, \Phi) + \theta_\alpha q_\alpha \boldsymbol{I} + \sum_{\beta \in C} \theta_\alpha p_{\alpha,\beta} \boldsymbol{R}(\beta, \Phi) - \lambda_\alpha \boldsymbol{I} - \theta_\alpha \boldsymbol{R}(\alpha, \Phi) \right]$$

$$= \Gamma \cdot \sum_{\alpha \in C} \left[ \left( \lambda_\alpha + \sum_{\beta \in C} \theta_\beta p_{\beta,\alpha} - \theta_\alpha \right) \boldsymbol{R}(\alpha, \Phi) + (\theta_\alpha q_\alpha - \lambda_\alpha) \boldsymbol{I} \right].$$

From this, we derive (5.2) to have $\Gamma \boldsymbol{A} = 0$, and it follows from (5.2) that $\sum_{\alpha \in C} \lambda_\alpha = \sum_{\alpha \in C} \theta_\alpha q_\alpha$.

Given hypothesis (3), we use properties (3) and (4) of Proposition 5.1 to show that

$$\Gamma \cdot \sum_{\alpha \in C} \lambda_\alpha \boldsymbol{R}(\alpha, \Phi) = \Gamma \cdot \sum_{\alpha \in C} \frac{\lambda_\alpha \mu_\alpha}{\theta_\alpha} \boldsymbol{M}(\alpha, \Phi)$$

$$= \kappa \cdot \sum_{\alpha \in C} \Gamma \cdot \boldsymbol{M}(\alpha, \Phi)$$

$$= \kappa \sum_{\alpha \in C} \Gamma \cdot \boldsymbol{M}(\alpha, \Psi)$$

$$= \Gamma \cdot \sum_{\alpha \in C} \lambda_\alpha \boldsymbol{R}(\alpha, \Psi).$$

This substitution then gives us a generator from an $\boldsymbol{M^C}/\boldsymbol{M^C}/\boldsymbol{1}//\boldsymbol{GD}$ queue where the insertion and deletion disciplines are equivalent to $\Psi$. From the previous arguments, it is clear that $\Gamma \boldsymbol{A} = 0$.

The above results show that we can easily characterize the steady-state distribution for a special subclass of $M^C/M^C/1//GD$ queues. They correspond to the quasi-reversible systems that are discussed in Kelly [2]. We now proceed to derive stochastic bounds that hold for *all* $M^C/M^C/1//GD$ systems. As we did for the $(M/M/1)^N$ network, we must construct an appropriate $S$ operator.

*Proposition* 5.3. *Define $S$ as a linear map from $\mathcal{F}(l_1(C))$ to $l_1$ where*

$$\left(\bigotimes_{i=1}^{n} e_{\alpha_i}\right)S = e_n$$

*we then have for all classes $\alpha$ and disciplines $\Phi$:*

(1) $R(\alpha, \Phi)S = SR$,
(2) $\sum_{\alpha \in C} L(\alpha, \Phi)S = SL$,
(3) $M(\alpha, \beta, \Phi)S = L(\alpha, \Phi)SR$.

*Proof.* It is sufficient to prove these statements for an arbitrary $\bigotimes_{i=1}^{n} e_{\alpha_i}$, for example

$$\left(\bigotimes_{i=1}^{n} e_{\alpha_i}\right)R(\alpha, \Phi)S$$

$$= \sum_{i=1}^{n+1} \phi_i(\alpha_1, \cdots, \alpha_n, \alpha)\left(\bigotimes_{j<i} e_{\alpha_j}\right) \otimes e_\alpha \otimes \left(\bigotimes_{j \geq i} e_{\alpha_j}\right)S$$

$$= \sum_{i=1}^{n+1} \phi_i(\alpha_1, \cdots, \alpha_n, \alpha)e_{n+1}$$

$$= e_{n+1}$$

$$= \left(\bigotimes_{i=1}^{n} e_{\alpha_i}\right)SR.$$

Similar proofs hold for Equations (2) and (3).

*Theorem* 5.4. *For a given $M^C/M^C/1//GD$ system, let $Q(t)$ be its queue-length process. Let $X(t)$ and $Y(t)$ be $M/M/1$ queue-length processes with both having arrival rate $\sum_{\alpha \in C} \lambda_\alpha$, and the service rate for each one is $\min_{\alpha \in C} \mu_\alpha q_\alpha$ and $\max_{\alpha \in C} \mu_\alpha q_\alpha$ respectively. If $Q(0) = X(0) = Y(0)$, then*

$$X(t) \geq_{st} Q(t) \geq_{st} Y(t)$$

*for all time $t \geq 0$.*

*Proof.* Let $\quad B_i = \lambda_i R + \mu_i L - \lambda_i I - \mu_i LR \quad$ where $\quad \lambda_1 = \lambda_2 = \sum_{\alpha \in C} \lambda_\alpha, \quad \mu_1 = \min_{\alpha \in C} \mu_\alpha q_\alpha$, and $\mu_2 = \max_{\alpha \in C} \mu_\alpha q_\alpha$. $X(t)$ and $Y(t)$ are monotone processes with $B_1$ the generator for $X(t)$, and $B_2$ the generator for $Y(t)$. By an argument similar to Theorem 4, it is sufficient to show that

$$SB_1K \geq ASK \geq SB_2K.$$

By the above proposition, and using the properties of $\mathbf{K}$, we get

$$
\begin{aligned}
\mathbf{ASK} &= \sum_{\alpha \in C} \left[ \lambda(\alpha)\mathbf{SR} + \mu_\alpha q_\alpha \mathbf{L}(\alpha, \Psi)\mathbf{S} \right. \\
&\quad \left. + \sum_{\beta \in C} \mu_\alpha p_{\alpha\beta}\mathbf{L}(\alpha, \Psi)\mathbf{SR} - \lambda_\alpha \mathbf{S} - \mu_\alpha \mathbf{L}(\alpha, \Psi)\mathbf{SR} \right] \mathbf{K} \\
&= \sum_{\alpha \in C} [\lambda_\alpha \mathbf{S}(\mathbf{R} - \mathbf{I}) + \mu_\alpha q_\alpha \mathbf{L}(\alpha, \Psi)\mathbf{S}(\mathbf{I} - \mathbf{R})]\mathbf{K} \\
&= \sum_{\alpha \in C} [\lambda_\alpha \mathbf{SR} - \mu_\alpha q_\alpha \mathbf{L}(\alpha, \Psi)\mathbf{SR}].
\end{aligned}
$$

Recall that $\mathbf{B}_i\mathbf{K} = \lambda_i\mathbf{R} - \mu_i\mathbf{LR}$ and $\sum_{\alpha \in C} \mathbf{L}(\alpha, \Psi)\mathbf{S} = \mathbf{SL}$. It is then a simple matter to show that $\mathbf{SB}_1\mathbf{K} \geqq \mathbf{ASK} \geqq \mathbf{SB}_2\mathbf{K}$.

Notice that we could have chosen initial distributions such that $X(0) \geqq_{st} Q(0) \geqq_{st} Y(0)$, and the theorem would still hold. We merely chose distributions that would give the closest fit. If we applied the above theorem to the previous example of an $\mathbf{M}/\mathbf{H}_2/\mathbf{1}$ queue, then it would have an $\mathbf{M}/\mathbf{M}/\mathbf{1}$ 'upper bound' with arrival rate $\lambda_\alpha + \lambda_\beta$ and service rate $\min(\mu_\alpha, \mu_\beta)$. The 'lower bound' would have the same arrival rate, but have service rate $\max(\mu_\alpha, \mu_\beta)$.

## 6. The $(\mathbf{M}^C/\mathbf{M}^C/\mathbf{1}//\mathbf{GD})^N$ network

The generators for the $(\mathbf{M}/\mathbf{M}/\mathbf{1})^N$ network (4.1) and the $\mathbf{M}^C/\mathbf{M}^C/\mathbf{1}//\mathbf{GD}$ queue (5.1) are very much alike in form. This illustrates the similarities between the notions of classes and nodes. In fact, they differ in only one way. There is no sense of order between two customers at different nodes for the $(\mathbf{M}/\mathbf{M}/\mathbf{1})^N$ network. On the other hand, the complexities of disciplines are introduced for the $\mathbf{M}^C/\mathbf{M}^C/\mathbf{1}//\mathbf{GD}$ queue to handle this notion of ordering. Algebraically, this is the difference between $E$ being a free non-abelian group on $C$ and being a free *abelian* group. If it is the latter, then we have a semigroup structure that makes $E$ isomorphic to $|C|$-tuples of non-negative integers. Setting $|C| = N$ gives us then the state space for the $(\mathbf{M}/\mathbf{M}/\mathbf{1})^N$ network.

We shall now proceed to construct the generator for the $(\mathbf{M}^C/\mathbf{M}^C/\mathbf{1}//\mathbf{GD})^N$ network, thereby defining the system. It is essentially an $(\mathbf{M}/\mathbf{M}/\mathbf{1})^N$ network modified so that each node is an $\mathbf{M}^C/\mathbf{M}^C/\mathbf{1}//\mathbf{GD}$ queue. Let $C$ be the set of classes for the *entire* network. We shall decompose $C$ into $N$ disjoint sets $D_1$ through $D_N$ or $C = \bigcup_{i=1}^N D_i$ with $D_i \cap D_j = \varnothing$ for $i \neq j$. Now let $E = \bigoplus_{i=1}^N S_{D_i}$, where we take the direct product of the $S_{D_i}$'s as semigroups. This is equivalent to generating a semigroup on $C$ with additional commutation relations. We

want $D_i$ to tag the classes that enter and exist only in the $i$th node. This is modelled by permitting any element of $D_i$ to commute with all of the elements of $C$ except the other members of $D_i$. If $\sigma_i$ is a string belonging to $S_{D_i}$, then $\sigma_1 \oplus \cdots \oplus \sigma_N$ is an element of $E$. This state space is then rich enough to describe an $N$-node network with various classes of customers served at each node.

It is clear from our previous construction that probability distributions on $E$ are represented by vectors belonging to the space

$$l_1(E) \cong l_1\left(\bigotimes_{i=1}^{N} S_{D_i}\right) \cong \bigotimes_{i=1}^{N} \mathscr{F}(l_1(D_i)).$$

Thus we have a tensor product of Fock spaces. This makes it easy to define operators on $l_1(E)$. Let $\boldsymbol{R}_i(\alpha, \Phi_i)$, with $\alpha$ belonging to $D_i$, be the operator that corresponds to inserting an $\alpha$-class customer into the $i$th node according to the discipline $\Phi_i$. We can construct such an $\boldsymbol{R}(\alpha, \Phi_i)$ to act on $\mathscr{F}(l_1(D_i))$. For $l_1(E) = \bigotimes_{i=1}^{N} \mathscr{F}(l_1(D_i))$, we define $\boldsymbol{R}_i(\alpha, \Phi_i)$ as

$$\boldsymbol{R}_i(\alpha, \Phi_i) = \boldsymbol{I} \otimes \cdots \otimes \boldsymbol{R}(\alpha, \Phi_i) \otimes \cdots \otimes \boldsymbol{I} \quad (i\text{th place}).$$

In a similar fashion, we define $\boldsymbol{L}_i(\alpha, \Psi_i)$ and $\boldsymbol{M}_i(\alpha, \beta, \Psi_i)$ for all $\alpha$ and $\beta$ belonging to $D_i$.

Just as for the $\boldsymbol{M^C/M^C/1//GD}$ queue, $\lambda_\alpha$ and $\mu_\alpha$ are, respectively, the arrival and service rates for the $\alpha$-class customers. Moreover, $\{p_{\alpha\beta}\}_{\alpha \in C, \beta \in C}$ is the $|C| \times |C|$ switching matrix that is substochastic so $q_\alpha = 1 - \sum_{\beta \in C} p_{\alpha\beta}$. If $\alpha$ and $\beta$ both belong to $D_i$, then $p_{\alpha\beta}$ describes the probability of an internal switch in node $i$, from class $\alpha$ to $\beta$. So when $\alpha$ is in $D_i$ and $\beta$ is in $D_j$ for $i \neq j$, then $p_{\alpha\beta}$ describes the probability of an external switch from node $i$ to node $j$ as well as a change of class. Finally, if $(\Phi_i, \Psi_i)$ is the insertion–deletion discipline for the $i$th node, then $\boldsymbol{A}$ equals

(6.1)
$$\boldsymbol{A} = \sum_{i=1}^{N} \sum_{\alpha \in D_i} \left[ \lambda_\alpha \boldsymbol{R}_i(\alpha, \Phi_i) + \mu_\alpha q_\alpha \boldsymbol{L}_i(\alpha, \Psi_i) + \sum_{\beta \in D_i} \mu_\alpha p_{\alpha\beta} \boldsymbol{M}_i(\alpha, \beta, \Psi_i) \right.$$
$$\left. + \sum_{j \neq i}^{N} \sum_{\beta \in D_i} \mu_\alpha p_{\alpha\beta} \boldsymbol{L}_i(\alpha, \Psi_i) \boldsymbol{R}_j(\beta, \Phi_j) - \lambda_\alpha \boldsymbol{I} - \mu_\alpha \boldsymbol{M}_i(\alpha, \Psi_i) \right].$$

We now consider the subclass of $(\boldsymbol{M^C/M^C/1//GD})^N$ networks which admit a steady distribution that is a tensor product of tensor geometric distributions. First, start with $\bigotimes_{i=1}^{N} \Gamma_i$ where each $\Gamma_i$ belongs to $\mathscr{F}(S_{D_i})$ and $(\bigotimes_{i=1}^{N} \Gamma_i)\boldsymbol{A} = 0$. By the same arguments as for the $(\boldsymbol{M/M/1})^N$ system, we have for each $\Gamma_i$

$$\Gamma_i \sum_{\alpha \in D_i} [\xi_\alpha \boldsymbol{R}(\alpha, \Phi_i) + \eta_\alpha \boldsymbol{L}(\alpha, \Psi_i) - \xi_\alpha \boldsymbol{I} - \eta_\alpha \boldsymbol{M}(\alpha, \Psi_i)] = 0$$

for some collection of positive scalars $\xi_\alpha$ and $\eta_\alpha$. So whenever $\bigotimes_{i=1}^{N} \Gamma_i$ is the steady-state distribution for an $(M^C/M^C/1//GD)^N$ system, then each $\Gamma_i$ is the steady-state distribution for some associated $M^C/M^C/1//GD$ system. This motivates the following theorem.

*Theorem* 6.1. *Let* $\Gamma_i$ *be a tensor geometric distribution belonging to* $\mathcal{F}(l_1(D_i))$ *generated by* $\mathbf{g}_i = \sum_{\alpha \in D_i} \dfrac{\theta_\alpha}{\mu_\alpha} \mathbf{e}_\alpha$. *Suppose that for each node i, one of the following conditions holds:*

(1) $\Phi_i$ *and* $\Psi_i$ *are abelian disciplines with* $\Phi_i \simeq \Psi_i$ *and* $p_{\alpha\beta} = 0$ *for all* $\alpha$ *and* $\beta$ *in* $D_i$.

(2) $\Phi_i$ *and* $\Psi_i$ *are load-dependent disciplines with* $\Phi_i \simeq \Psi_i$.

*If the steady-state distribution exists, then* $(\bigotimes_{i=1}^{N} \Gamma_i)\mathbf{A} = 0$, *where the* $\theta_\alpha$'s *solve*

$$(6.2) \qquad \theta_\alpha = \lambda_\alpha + \sum_{\beta \in C} \theta_\beta p_{\beta\alpha}$$

*and* $\sum_{\alpha \in D_i} \dfrac{\theta_\alpha}{\mu_\alpha} < 1$ *for each i.*

*Proof.* Using the results of Proposition 5, hypothesis (1) or (2), and the properties of tensor products gives us

$$\left(\bigotimes_{i=1}^{N} \Gamma_i\right)\mathbf{A} = \left(\bigotimes_{i=1}^{N} \Gamma_i\right) \sum_{i=1}^{N} \sum_{\alpha \in D_i} \left[ \lambda_\alpha \mathbf{R}_i(\alpha, \Phi_i) + \theta_\alpha q_\alpha \mathbf{I} \right.$$

$$\left. + \sum_{j=1}^{N} \sum_{\beta \in D_i} \theta_\alpha p_{\alpha\beta} \mathbf{R}_j(\beta, \Phi_j) - \lambda_\alpha \mathbf{I} - \theta_\alpha \mathbf{R}_i(\alpha, \Phi_i) \right]$$

$$= \left(\bigotimes_{i=1}^{N} \Gamma_i\right) \sum_{i=1}^{N} \sum_{\alpha \in D_i} \left[ \left( \lambda_\alpha + \sum_{\beta \in C} \theta_\beta p_{\beta\alpha} - \theta_\alpha \right) \mathbf{R}_i(\alpha, \Phi_i) + (\theta_\alpha q_\alpha - \lambda_\alpha)\mathbf{I} \right]$$

and as in the proof of Theorem 5.2, to have $\left(\bigotimes_{i=1}^{N} \Gamma_i\right)\mathbf{A} = 0$ requires that the $\theta_\alpha$'s satisfy (6.2). From this follows $\sum_{\alpha \in C} \lambda_\alpha = \sum_{\alpha \in C} \theta_\alpha q_\alpha$.

*Corollary* 6.2. *If in addition to condition* (1) *or* (2) *above, some nodes satisfy*

(3) $\Phi_i$ *and* $\Psi_i$ *are load-dependent with* $\Phi_i \neq \Psi_i$.

*Then, we still have* $(\bigotimes_{i=1}^{N} \Gamma_i)\mathbf{A} = 0$, *if whenever the ith node satisfies condition* (3), *there is a constant* $\kappa_i$ *such that for all* $\alpha$ *in* $D_i$, $\mu_\alpha = \dfrac{\theta_\alpha}{\theta_\alpha - \sum\limits_{\beta \in D_i} \theta_\beta p_{\beta\alpha}} \cdot \kappa_i$, *and the* $\theta_\alpha$'s *still solve* (6.2).

*Proof.* It follows from Proposition 5.1 that for all $j$ the quantity $(\bigotimes_{i=1}^{N} \Gamma_i)$ $\cdot \sum_{\alpha \in D_i} M_j(\alpha, \Psi_i)$ is independent of $\Psi_j$. If the $i$th node satisfies condition (3), then we can replace $R_i(\alpha, \Psi_i)$ by $R_i(\alpha, \Phi_i)$ as we did the opposite in the proof of Theorem 5.2. The rest of the proof would follow the same steps as above. Knowing the $\theta_\alpha$'s in advance, we would then get 0.

The special class of $(M^C/M^C/1//GD)^N$ networks characterized above, encompass the product-form networks as discussed in Kelly [2].

We now proceed to derive the main theorem of this paper, which derives stochastic bounds for *any* $(M^C/M^C/1//GD)^N$ network, thereby estimating its transient behavior by that of $M/M/1$ queues. We first define two aggregation operators $S_*$ and $S$. $S_*$ maps $\bigotimes_{i=1}^{N} \mathcal{F}(l_1(D_i))$ to $l_1^{(N)}$ where

$$(e_{\sigma_1} \otimes \cdots \otimes e_{\sigma_N})S_* = e_{|\sigma_1|} \otimes \cdots \otimes e_{|\sigma_N|}$$

with $\sigma_i$ being an arbitrary string in $S_{D_i}$ and $|\sigma_i|$ is its length. $S$ on the other hand, maps $\bigotimes_{i=1}^{N} \mathcal{F}(l_1(D_i))$ into $l_1$ where

$$(e_{\sigma_1} \otimes \cdots \otimes e_{\sigma_N})S = e_{|\sigma_1| + \cdots + |\sigma_N|}$$

we now list the key properties of $S_*$ and $S$. The method of proof is similar to that for Proposition 5.3, so we shall omit it.

*Proposition 6.3. For all classes $\alpha$ and disciplines $\Phi_i$, we have for $S_*$*
(1) $R_i(\alpha, \Phi_i)S_* = S_*R_i$,
(2) $\sum_{\alpha \in D_i} L_i(\alpha, \Phi_i)S_* = S_*L_i$,
(3) $M_i(\alpha, \beta, \Phi_i)S_* = L_i(\alpha, \Phi_i)S_*R_i$,
*where $R_i$ and $L_i$ are defined on $l_1^{(N)}$. Similarly, for $S$ we have*
(4) $R_i(\alpha, \Phi_i)S = SR$,
(5) $\sum_{\alpha \in D_i} L_i(\alpha, \Phi_i)S \leq SL$,
(6) $M_i(\alpha, \beta, \Phi_i)S = L_i(\alpha, \Phi_i)SR$.

Before we prove our main results, we introduce the following set of variations on the $K$ operator to act on $l_1^{(N)}$:
(1) $K_i = I \otimes \cdots \otimes K \otimes \cdots \otimes I$ ($i$ th place).
(2) $K_I = \prod_{i \in I} K_i$ and $K_{(I)} = \prod_{i \notin I} K_i$ where $I$ is any subset of $\{1, \cdots, N\}$,
(3) $K_* = \prod_{i=1}^{N} K_i$.
In Massey [6], the author used $K_*$ as the $(M/M/1)^N$ network analogue to the $K$ operator for the $M/M/1$ queue. $K_*$ inherits the following algebraic properties from $K$:

$$R_i K_* = K_* + K_{(i)}R_i \quad \text{and} \quad L_i K_* = K_* - K_{(i)}.$$

Similar results hold for the $K_I$'s.

**Theorem 6.4.** *Let $X_1(t), \cdots, X_N(t)$ be a collection of independent **M/M/1** queue-length processes where $X_i(t)$ has arrival rate $\lambda_i$ and service rate $\mu_i$, where in terms of the given $(\mathbf{M^C/M^C/1//GD})^N$ network, we have*

$$\lambda_i = \sum_{\alpha \in D_i} \left( \lambda_\alpha + \sum_{j \neq i} \sup_{\beta \in D_j} \mu_\beta p_{\beta\alpha} \right)$$

$$\mu_i = \inf_{\alpha \in D_i} \mu_\alpha \left( 1 - \sum_{\beta \in D_i} p_{\alpha\beta} \right).$$

*If $X_i(0) = Q_i(0)$ for all i, then*

$$\Pr\{Q_1(t) \geq n_1, \cdots, Q_N(t) \geq n_N\} \leq \prod_{i=1}^{N} \Pr\{X_i(t) \geq n_i\}$$

*for all $t > 0$, and all non-negative integers $n_1, \cdots, n_N$.*

**Proof.** Let $\mathbf{B}_i = \lambda_i \mathbf{R}_i + \mu_i \mathbf{L}_i - \lambda_i \mathbf{I} - \mu_i \mathbf{L}_i \mathbf{R}_i$. Since the $X_i(t)$'s are independent, then $(X_1(t), \cdots, X_N(t))$ is a Markov process and its generator is $\sum_{i=1}^{N} \mathbf{B}_i$, which acts on $l_1^{(N)}$. If $\mathbf{A}$ is the generator for the $(\mathbf{M^C/M^C/1//GD})^N$ network, as defined in (6.1), then $\mathbf{p}(0) \exp(t\mathbf{A})\mathbf{S}_*$ encodes probabilities of the form $\Pr\{Q_1(t) = n_1, \ldots, Q_n(t) = n_N\}$. The $Q_i(t)$'s here are the queue-length process but in general $(Q_1(t), \ldots, Q_N(t))$ is *not* a Markov process. Regardless, $\mathbf{p}(0) \exp(t\mathbf{A})\mathbf{S}_*\mathbf{K}_*$ encodes probabilities of the form $\Pr\{Q_1(t) \geq n_1, \ldots, Q_N(t) \geq n_N\}$.

Given this, we need only show that

$$\exp(t\mathbf{A})\mathbf{S}_*\mathbf{K}_* \leq \mathbf{S}_* \exp\left( t \cdot \sum_{i=1}^{N} \mathbf{B}_i \right)\mathbf{K}_*.$$

By previous arguments of this type, we need only show that

$$\mathbf{A}\mathbf{S}_*\mathbf{K}_* \leq \mathbf{S}_* \cdot \sum_{i=1}^{N} \mathbf{B}_i \cdot \mathbf{K}_*$$

*provided* that $\mathbf{K}_*^{-1} \exp(t \cdot \sum_{i=1}^{N} \mathbf{B}_i)\mathbf{K}_* \geq 0$. The $\mathbf{B}_i$'s are defined so that they commute, hence

$$\exp\left( t \cdot \sum_{i=1}^{N} \mathbf{B}_i \right) = \prod_{i=1}^{N} \exp(t\mathbf{B}_i)$$

and so

$$\mathbf{K}_*^{-1} \exp\left( t \cdot \sum_{i=1}^{N} \mathbf{B}_i \right)\mathbf{K}_* = \prod_{i=1}^{N} \mathbf{K}_i^{-1} \exp(t\mathbf{B}_i)\mathbf{K}_i$$

$$= \bigotimes_{i=1}^{N} \mathbf{K}^{-1}\exp(t(\lambda_i \mathbf{R} + \mu_i \mathbf{L} - \lambda_i \mathbf{I} - \mu_i \mathbf{L}\mathbf{R}))\mathbf{K}.$$

As we showed in Section 3, each $K^{-1} \exp(t(\lambda_i R + \mu_i L - \lambda_i I - \mu_i LR))K$ is positive, so their product is positive also.

We can now proceed to show that

$$AS_* K_* \leqq S_*\left(\sum_{i=1}^{N} B_i\right)K_*.$$

Using the algebraic properties of $S_*$ and $K_*$ gives

$$
\begin{aligned}
AS_* K_* &= \sum_{i=1}^{N} \sum_{\alpha \in D_i} \left[ \lambda_\alpha S_* R_i + \mu_\alpha q_\alpha L_i(\alpha, \Psi_i) S_* + \sum_{\beta \in D_i} \mu_\alpha p_{\alpha\beta} L_i(\alpha, \Psi_i) S_* R_i \right. \\
&\quad \left. + \sum_{\substack{j \neq i \\ \beta \in D_i}}^{N} \mu_\alpha p_{\alpha\beta} L_i(\alpha, \Psi_i) S_* R_j - \lambda_\alpha S_* - \mu_\alpha L_i(\alpha, \Psi_i) S_* R_i \right] K_* \\
&= \sum_{i=1}^{N} \sum_{\alpha \in D_i} \left[ \lambda_\alpha S_*(K_{(i)} R_i + K_*) + \mu_\alpha q_\alpha L_i(\alpha, \Psi_i) S_* K_* \right. \\
&\quad + \sum_{\beta \in D_i} \mu_\alpha p_{\alpha\beta} L_i(\alpha, \Psi_i) S_*(K_{(i)} R_i + K_*) \\
&\quad + \sum_{\substack{j \neq i \\ \beta \in D_i}}^{N} \mu_\alpha p_{\alpha\beta} L_i(\alpha, \Psi_i) S_*(K_{(j)} R_j + K_*) \\
&\quad \left. - \lambda_\alpha S_* K_* - \mu_\alpha L_i(\alpha, \Psi_i) S_*(K_{(i)} R_i + K_*) \right] \\
&= \sum_{i=1}^{N} \sum_{\alpha \in D_i} \left[ \lambda_\alpha S_* K_{(i)} R_i + \mu_\alpha \left( \sum_{\alpha \in D_i} p_{\alpha\beta} \right) L_i(\alpha, \Psi_i) S_* K_{(i)} R_i \right. \\
&\quad \left. + \sum_{\substack{j \neq i \\ \beta \in D_i}}^{N} \mu_\alpha p_{\alpha\beta} L_i(\alpha, \Psi_i) S_* K_{(j)} R_j - \mu_\alpha L_i(\alpha, \Psi_i) S_* K_{(i)} R_i \right] \\
&\leqq \sum_{i=1}^{N} \left[ \left( \sum_{\alpha \in D_i} \lambda_\alpha \right) S_* K_{(i)} R_i - \inf_{\alpha \in D_i} \mu_\alpha \left( 1 - \sum_{\beta \in D_i} p_{\alpha\beta} \right) S_* K_{(i)} L_i R_i \right. \\
&\quad \left. + \sum_{\substack{j \neq i \\ \beta \in D_i}}^{N} \left( \sup_{\alpha \in D_i} \mu_\alpha p_{\alpha\beta} \right) S_*(K_{(j)} R_j - K_{(i,j)} R_j) \right] \\
&\leqq \sum_{i=1}^{N} \left[ \lambda_i S_* K_{(i)} R_i - \mu_i S_* K_{(i)} L_i R_i - \sum_{\substack{j \neq i \\ \beta \in D_i}}^{N} \left( \sup_{\alpha \in D_i} \mu_\alpha p_{\alpha\beta} \right) S_* K_{(i,j)} R_j \right] \\
&\leqq S_* \cdot \sum_{i=1}^{N} B_i \cdot K_*
\end{aligned}
$$

and this finishes the proof.

*Theorem 6.5. Let $Y(t)$ be an **M/M/1** queue-length process with arrival rate $\lambda$ and service rate $\mu$ where*

$$\lambda = \sum_{i=1}^{N} \sum_{\alpha \in D_i} \lambda_\alpha \qquad \mu = \sum_{i=1}^{N} \sup_{\alpha \in D_i} \mu_\alpha q_\alpha.$$

*If* $Y(0) = \sum_{i=1}^{N} Q_i(0)$, *then*

$$\Pr\left\{\sum_{i=1}^{N} Q_i(t) \geqq n\right\} \geqq \Pr\{Y(t) \geqq n\}$$

*for all* $t > 0$ *and all non-negative integers* $n$.

*Proof.* Let $\boldsymbol{B} = \lambda\boldsymbol{R} + \mu\boldsymbol{L} + \lambda\boldsymbol{L} - \boldsymbol{LR}$. $\boldsymbol{B}$ represents a monotone process, so we need only show that $\boldsymbol{ASK} \geqq \boldsymbol{SBK}$.

$$
\begin{aligned}
\boldsymbol{ASK} &= \sum_{i=1}^{N} \sum_{\alpha \in D_i} \left[ \lambda_\alpha \boldsymbol{SR} + \mu_\alpha q_\alpha \boldsymbol{L}_i(\alpha, \Psi_i)\boldsymbol{S} + \sum_{j=1}^{N} \sum_{\beta \in D_i} \mu_\alpha p_{\alpha\beta} \boldsymbol{L}_i(\alpha, \Psi_i)\boldsymbol{SR} \right.\\
&\quad \left. - \lambda_\alpha \boldsymbol{S} - \mu_\alpha \boldsymbol{L}_i(\alpha, \Psi_i)\boldsymbol{SR} \right]\boldsymbol{K} \\
&= \sum_{i=1}^{N} \sum_{\alpha \in D_i} [\lambda_\alpha \boldsymbol{SR} - \mu_\alpha q_\alpha \boldsymbol{L}_i(\alpha, \Psi_i)\boldsymbol{SR}] \\
&\geqq \sum_{i=1}^{N} \left[ \sum_{\alpha \in D_i} \lambda_\alpha \boldsymbol{SR} - \sup_{\alpha \in D_i} \mu_\alpha q_\alpha \boldsymbol{SLR} \right] \\
&\geqq \lambda\boldsymbol{SR} - \mu\boldsymbol{SLR} \\
&\geqq \boldsymbol{SBK}
\end{aligned}
$$

and we are done.

## Acknowledgements

## References

[1] FAYOLLE, G., MITRANI, I. AND IASNOGORODSKI, R. (1980) Sharing a processor among many job classes. *J. Assoc. Comput. Mach.* **27**, 519–532.

[2] KELLY, F. P. (1976) Networks of queues. *Adv. Appl. Prob.* **8**, 416–432.

[3] KIRSTEIN, B. M., FRANKEN, P. AND STOYAN, D. (1977) Comparability and monotonicity of Markov processes. *Theory Prob. Appl.* **22**, 40–52.

[4] MASSEY, W. A. (1984) Asymptotic analysis of $M(t)/M(t)/1$: The time-dependent $M/M/1$ queue. *Math. Operat. Res.* To appear.

[5] MASSEY, W. A. (1984) An operator analytic approach to the Jackson network. *J. Appl. Prob.* **21** (2).