

POINT PROCESSES IN FAST JACKSON NETWORKS¹

BY JAMES B. MARTIN

INRIA and Ecole Normale Supérieure, Paris

We consider a Jackson-type network, each of whose nodes contains N identical channels with a single server. Upon arriving at a node, a task selects m of the channels at random and joins the shortest of the m queues observed. We fix a collection of channels in the network, and analyze how the queue-length processes at these channels vary as $N \rightarrow \infty$. If the initial conditions converge suitably, the distribution of these processes converges in local variation distance to a limit under which each channel evolves independently. We discuss the limiting processes which arise, and in particular we investigate the point processes of arrivals and departures at a channel when the networks are in equilibrium, for various values of the system parameters.

1. Introduction. The class of *Jackson networks* was introduced in [4] and [5] and has since been one of the most widely studied in queueing network theory. The basic model consists of a network of J nodes; each node j , $1 \leq j \leq J$, has an infinite buffer and a single server with service rate μ_j ; tasks arrive at node j from outside the network as a Poisson process of rate λ_j , and when a task completes its service at node j , it immediately joins the queue at node k with probability p_{jk} , $1 \leq k \leq J$, and leaves the network with probability $p_j^* = 1 - \sum_k p_{jk}$. All service times, external arrival times and routing decisions are independent.

One of the attractions of this model is the appealing *product form* of the stationary distribution for the network; in equilibrium (if the network is not overloaded), the lengths of the queues at the various nodes in the network at a given point in time are independent, and are each distributed as for an equilibrium $M/M/1$ queue with appropriately chosen arrival and service rates. However, this simplicity of description for single points in time does not extend to the pathwise behavior of the network. As specified, the arrivals from outside the network to the various nodes form a family of independent Poisson processes, and so by a reversibility argument (see, e.g., [7]) one can show that in equilibrium the same is true of the processes of departures leaving the network from the various nodes. However, the process of all arrivals, or of all departures, at a given node is not in general easy to describe, and of course the queue-length processes at different nodes of the network are not independent. Similarly, the joint distribution of the waiting times experienced

Received January 1999; revised October 2000.

¹Supported in part by EPSRC studentship no. 95007341 and by European Union TMR grant ALAPEDES (RB-FMRX-CT-96-0074).

AMS 2000 *subject classifications*. 90B15, 60G55.

Key words and phrases. Queueing network, point process, Jackson network.

by a particular task at the various nodes it visits in the network is not easy to obtain.

We consider a modification of the Jackson network model, in which each node j of the network contains N channels, each with an infinite buffer and a single server with service rate μ_j . External arrivals arrive at node j at rate $N\lambda_j$, and routing between nodes is as before, according to the matrix $\mathbf{P} = (p_{jk})$. Upon arrival at a node (either from outside the network or after being served at the same or another node), each task now inspects m channels chosen uniformly at random from the N available (with replacement, though this is unimportant for large N), and joins the shortest of the m queues observed (breaking ties at random). The behavior of the network is therefore specified by the parameters J , N , $\underline{\lambda}$, $\underline{\mu}$, \mathbf{P} and m .

Systems corresponding to a single node in this model were studied by Vvedenskaya, Dobrushin and Karpelevich in [9]. The network model just described was introduced in [8]. There it is shown that, as $N \rightarrow \infty$, the evolution of the system may be asymptotically represented by the solution of a countably infinite system of ODEs; under a standard nonoverload condition on the parameters $\underline{\lambda}$, $\underline{\mu}$ and \mathbf{P} , the system has an invariant distribution π_N for each N , and, as $N \rightarrow \infty$, π_N converges to a limiting invariant distribution which is concentrated at a single point, corresponding to the fixed point of the system of ODEs. Under this limiting distribution, for $m > 1$, the tail of the distribution of queue lengths decays superexponentially rather than exponentially as in the case of standard Jackson networks; hence the term “Fast Jackson networks.”

In this paper we again let $N \rightarrow \infty$, but now consider how the paths of the queue length processes at individual channels behave as the size of the network grows. We show that, if the initial state of the network converges suitably, the distribution of the queue length processes at a fixed collection of channels at the same or different nodes converges in “local variation distance” as $N \rightarrow \infty$, and that under the limit the component processes are independent. We describe the limiting processes which arise, and analyze them in particular in the case where the networks are positive recurrent and are started in equilibrium. Then for $m = 1$ (in which case the networks are standard Jackson networks for finite N), the limiting point processes of arrivals and departures are Poisson processes, and we examine how they and the relationship between them change as m increases. For networks in equilibrium, another interpretation of the decoupling which occurs in the limit is that a typical task, given its route through the network, experiences a sequence of independent waiting times.

In the next section we introduce notation and restate results from [8] which we will use. In Section 3 the main theorem is proved, using results of Kabanov and Liptser from [6] which relate convergence in variation distance of multivariate point processes to the convergence of their compensators. In Section 4 we interpret this result for the case of networks in equilibrium, and analyze and illustrate the particular point processes that arise for various different values of the network parameters. Finally in Section 5 we discuss possible

extensions of the results; in particular we compare our approach with that of Brown and Pollett [2], who investigate how the distance of arrival processes in a standard Jackson network from appropriate Poisson processes varies as the number of nodes in the network is increased.

2. Preliminaries. The state of a network as described above with N channels at each node may be described by a vector $\mathbf{r} = \{r_j(n), 1 \leq j \leq J, n \in \mathbb{Z}_+\}$, (here and below \mathbb{Z}_+ is the set of nonnegative integers), where $r_j(n) = N^{-1} \sum_{n' \geq n} M_j(n')$ and $M_j(n')$ is the number of channels at node j whose queue length (including the customer in service) is n' . Hence $r_j(n)$ is the proportion of channels at node j whose queue length is at least n . The process $\mathbf{r}(t) = \{r_j(n, t), n \in \mathbb{Z}_+, t \geq 0\}$, describing the state of the network at times $t \geq 0$, is easily seen to be a Markov process for each N , with state space \mathcal{U}_N^J where

$$(2.1) \quad \mathcal{U}_N = \left\{ \mathbf{g} = (g(n), n \in \mathbb{Z}_+): g(0) = 1, g(n) \geq g(n+1) \geq 0, \right. \\ \left. Ng(n) \in \mathbb{N}, \forall n, \text{ and } g(n) = 0 \text{ for sufficiently large } n \right\}.$$

Since we wish to let $N \rightarrow \infty$, we will also consider the limiting space \mathcal{U}^J , where

$$\mathcal{U} = \left\{ \mathbf{g} = (g(n), n \in \mathbb{Z}_+): g(0) = 1, g(n) \geq g(n+1) \geq 0, \forall n, \right. \\ \left. \text{and } \sum_{n=0}^{\infty} g(n) < \infty \right\}.$$

Then $\mathcal{U}_N^J \subset \mathcal{U}^J$ for all N . Following [9], we define the metric

$$(2.2) \quad d(\mathbf{u}, \mathbf{u}') = \sup_{1 \leq j \leq J} \sup_{n \geq 1} \frac{|u_j(n) - u'_j(n)|}{n}$$

on the spaces \mathcal{U}^J and \mathcal{U}_N^J .

We will consider the following infinite system of nonlinear differential equations for $\mathbf{u}(t) = \{u_j(n, t), 1 \leq j \leq J, n \in \mathbb{Z}_+, t \geq 0\}$, with initial condition $\mathbf{g} \in \mathcal{U}^J$:

$$(2.3) \quad \mathbf{u}(0) = \mathbf{g},$$

$$(2.4) \quad \dot{\mathbf{u}}(t) = \mathbf{h}(\mathbf{u}(t)),$$

where, for all j ,

$$(2.5) \quad h_j(0, \mathbf{u}) = 0,$$

$$(2.6) \quad h_j(n, \mathbf{u}) = \left[\lambda_j + \sum_{1 \leq k \leq J} \mu_k p_{kj} u_k(1) \right] [u_j(n-1)^m - u_j(n)^m] \\ - \mu_j [u_j(n) - u_j(n+1)] \quad \text{for all } n \geq 1.$$

The following result, proved in [8], then describes how the solution of this system asymptotically represents the behavior of the network.

THEOREM 2.1. (i) *If $\underline{\mathbf{g}} \in \mathcal{U}^J$, the system (2.3)–(2.6) has a unique solution $\underline{\mathbf{u}}(t, \underline{\mathbf{g}})$, $t \geq 0$ in \mathcal{U}^J .*

(ii) *For any continuous function $f: \mathcal{U}^J \rightarrow \mathbb{R}$ and $t \geq 0$,*

$$\lim_{N \rightarrow \infty} \sup_{\underline{\mathbf{g}} \in \mathcal{U}_N^J} |\mathbb{E}_N[f(\mathbf{r}(s)) | \mathbf{r}(0) = \underline{\mathbf{g}}] - f(\underline{\mathbf{u}}(s, \underline{\mathbf{g}}))| = 0,$$

uniformly in $s \in [0, t]$, where \mathbb{E}_N denotes the expectation under the dynamics of the network with N channels at each node.

3. Convergence as $N \rightarrow \infty$. We fix the parameters $J, \underline{\lambda}, \underline{\mu}, \mathbf{P}$ and m , and consider a sequence of networks indexed by N , with the N th network having N channels at each node. We fix a set of K tagged channels among the various nodes; formally, we fix a function i from $\{1, 2, \dots, K\}$ to $\{1, 2, \dots, J\}$ such that, in each network, the k th tagged channel belongs to node $i(k)$, $1 \leq k \leq K$. We will analyze the behavior of the process of queue lengths at the tagged channels as $N \rightarrow \infty$.

We will describe the evolution of the network by constructing, for each N , a process $(\mathbf{r}(t), \mathbf{x}(t))$, $t \geq 0$, with state space $\mathcal{U}^J \times \mathbb{Z}_+^K$. Here $x_k(t)$ will be the length of the queue in the k th tagged channel at time t , and $r_j(n, t)$ will be the proportion of channels (including tagged channels) at node j whose queue length is at least n at time t . Thus we do not distinguish between different untagged channels at the same node. The topology used on $\mathcal{U}^J \times \mathbb{Z}_+^K$ is the product of the topology induced by the metric (2.2) on \mathcal{U}^J and the discrete topology on \mathbb{Z}_+^K .

Let Ω be the space of paths $[0, \infty) \rightarrow \mathcal{U}^J \times \mathbb{Z}_+^K$ which are right continuous with left limits, [representing the paths of $(\mathbf{r}(t), \mathbf{x}(t))$]. For each $t \geq 0$ we define the functions $\mathbf{r}(t)$ and $\mathbf{x}(t)$ on Ω by setting $(\mathbf{r}(t), \mathbf{x}(t))(w) = w(t)$. Define the σ -algebra \mathcal{S} on Ω by $\mathcal{S} = \sigma(\mathbf{r}(s), \mathbf{x}(s), s \geq 0)$.

For each $N \geq K$, let ψ_N be a distribution on $\mathcal{U}^J \times \mathbb{Z}_+^K$, representing the distribution of $(\mathbf{r}(0), \mathbf{x}(0))$ for the N th network. This, together with the dynamics of the N th network described earlier, yields a probability measure $\bar{\mathbb{P}}_N$ on the measurable space (Ω, \mathcal{S}) describing the behavior of the N th network. We write \mathbb{E}_N for the expectation with respect to $\bar{\mathbb{P}}_N$.

We define the filtration $G = \{\mathcal{G}_t\}_{t \geq 0}$ on \mathcal{S} by $\mathcal{G}_t = \sigma(\mathbf{r}(s), \mathbf{x}(s), 0 \leq s \leq t)$. For each N , $(\mathbf{r}(t), \mathbf{x}(t), t \geq 0)$ is a stochastic process defined on $(\Omega, \mathcal{S}, \bar{\mathbb{P}}_N)$ and adapted to the filtration G . Note that $\bar{\mathbb{P}}_N(\mathbf{r}(t) \in \mathcal{U}_N^J \forall t) = 1$. Ultimately we will be particularly interested in the smaller filtration $F = \{\mathcal{F}_t\}_{t \geq 0}$, where $\mathcal{F}_t = \sigma(\mathbf{x}(s), 0 \leq s \leq t)$, and $\mathcal{F} = \bigvee_t \mathcal{F}_t$. Clearly $\mathcal{F}_t \subset \mathcal{G}_t \forall t$ and $\mathcal{F} \subset \mathcal{S}$, and the process $\mathbf{x}(t)$ is adapted to the filtration F ; note also that since, for all $\omega \in \Omega$, the path $\mathbf{x}(t)(w)$ is right continuous with respect to the discrete

topology, the filtration F is itself right continuous. Let \mathbb{P}_N be the restriction of $\bar{\mathbb{P}}_N$ to (Ω, \mathcal{F}) .

For two measures \mathbb{P} and \mathbb{P}' on (Ω, \mathcal{F}) , we will write $\text{Var}_t(\mathbb{P}, \mathbb{P}')$ for the variation distance between \mathbb{P} and \mathbb{P}' restricted to \mathcal{F}_t ,

$$\text{Var}_t(\mathbb{P}, \mathbb{P}') = \sup_{A \in \mathcal{F}_t} |\mathbb{P}(A) - \mathbb{P}'(A)|.$$

The following theorem states that, if the initial conditions converge suitably, then the processes \mathbf{x} governed by \mathbb{P}_N converge in this local variation distance for each t .

THEOREM 3.1. *Suppose that $\psi_N \rightarrow \psi$ weakly, where ψ is a distribution on $\mathcal{W}^J \times \mathbb{Z}_+^K$ under which the marginal distribution of $\underline{\mathbf{x}}(0)$ on \mathcal{W}^J is concentrated at a single point. Then:*

- (i) *There exists a probability measure \mathbb{P} on (Ω, \mathcal{F}) such that for all $t \geq 0$, $\text{Var}_t(\mathbb{P}_N, \mathbb{P}) \rightarrow 0$ as $N \rightarrow \infty$.*
- (ii) *Under the limiting measure \mathbb{P} , $\{\mathbf{x}(t), t \geq 0\}$ is a Markov process (not in general time-homogeneous) and if $x_1(0), \dots, x_K(0)$ are independent under ψ , then the component processes $\{x_k(t), t \geq 0\}$, $1 \leq k \leq K$, are independent under \mathbb{P} .*

PROOF. Given the initial state $\mathbf{x}(0)$, we may represent the process $\mathbf{x}(t)$ by a multivariate point process $\xi = \{\xi_{k,l}(t), t \geq 0, k, l \in \{0, 1, \dots, K\}, (l, k) \neq (0, 0)\}$ with $K^2 + 2K$ components. Each component is an increasing integer-valued process, which has value 0 at time 0, and whose value at time t is the number of *points* which occur in that component during the time interval $(0, t]$. For $1 \leq k \leq K$, let the points of the component process $\xi_{0,k}(t)$ record times of arrivals at the k th tagged channel which do not originate from a tagged channel (so are external arrivals or tasks transferring from an untagged channel at the same or another node), and let those of $\xi_{k,0}(t)$ record times of departures from the k th tagged channel which do not proceed to another tagged channel. Finally, for $1 \leq k, l \leq K$, let the points of $\xi_{k,l}(t)$ record times of transfers from the k th tagged channel to the l th (which is the same if $k = l$). So, for example, the total number of departures from the k th tagged channel during the time interval $(0, t]$ is $\sum_{l=0}^K \xi_{k,l}(t)$. Note that, $\bar{\mathbb{P}}_N$ -a.s. for all N , no two points occur simultaneously, in the same or different components.

For each N , we associate with the point process ξ and the filtered probability space $(\Omega, \mathbb{P}_N, \mathcal{F}, F)$ to which it is adapted the multivariate *compensator* $\mathbf{B}^{(N)}$ (indexed in the same way as ξ) which is the (unique) F -previsible process such that $\mathbf{B}^{(N)}(0) = \mathbf{0}$ and each component of $\xi - \mathbf{B}^{(N)}$ is a (\mathbb{P}_N, F) -martingale. If the limit

$$(3.1) \quad \boldsymbol{\beta}^{(N)}(t) = \lim_{h \downarrow 0} h^{-1} \mathbb{E}_N(\xi(t+h) - \xi(t) | \mathcal{F}_t)$$

exists a.s. for all $t \geq 0$, then $\boldsymbol{\beta}^{(N)}$ is called the *conditional intensity process* for $\boldsymbol{\xi}$ (with respect to F), and we have

$$\mathbf{B}^{(N)}(t) = \int_0^t \boldsymbol{\beta}^{(N)}(s) ds \quad \text{a.s.},$$

for all t . (All such relations are to be understood componentwise). For further details, see, for example, Section 13.2 of [3].

To show the existence of $\boldsymbol{\beta}^{(N)}$ for all N , and to demonstrate the convergence as $N \rightarrow \infty$, we will additionally consider the conditional intensity of the point process $\boldsymbol{\xi}$ with respect to $(\Omega, \overline{\mathbb{P}}_N, \mathcal{G}, G)$, denoting this by $\boldsymbol{\alpha}^{(N)}$. Analogously to (3.1), we have

$$(3.2) \quad \boldsymbol{\alpha}^{(N)}(t) = \lim_{h \downarrow 0} h^{-1} \mathbb{E}_N(\boldsymbol{\xi}(t+h) - \boldsymbol{\xi}(t) | \mathcal{G}_t).$$

Since, for each N , $(\mathbf{r}(s), \mathbf{x}(s))$ is a countable state-space Markov process under $\overline{\mathbb{P}}_N$, with $\mathcal{G}_t = \sigma(\{\mathbf{r}(s), \mathbf{x}(s)\}, 0 \leq s \leq t)$, the limit (3.2) exists a.s. for all t , and corresponds to a vector of certain instantaneous transition rates from the state $(\mathbf{r}(t), \mathbf{x}(t))$; the transitions concerned are those representing a departure or an arrival at a tagged channel, giving rise to a point in the process $\boldsymbol{\xi}$.

Now for any state of the network, the instantaneous departure rate from any tagged channel is no greater than $\max_i \mu_i$, and the instantaneous arrival rate to any tagged channel (from outside the network and from other nodes in the network) is no greater than

$$m \max_i \left\{ \lambda_i + \sum_j \mu_j p_{ji} \right\}$$

(which is m/N times the maximal arrival rate at any single node). Hence we have

$$h^{-1} \mathbb{E}_N(\xi_{k,l}(t+h) - \xi_{k,l}(t) | \mathcal{G}_t) \leq \max_i \mu_i + m \max_i \left(\lambda_i + \sum_j \mu_j p_{ji} \right)$$

a.s. for all N , all k, l , all t and all $h > 0$.

Then

$$(3.3) \quad \begin{aligned} \boldsymbol{\beta}^{(N)}(t) &= \lim_{h \downarrow 0} h^{-1} \mathbb{E}_N(\boldsymbol{\xi}(t+h) - \boldsymbol{\xi}(t) | \mathcal{F}_t) \\ &= \lim_{h \downarrow 0} h^{-1} \mathbb{E}_N(\mathbb{E}_N(\boldsymbol{\xi}(t+h) - \boldsymbol{\xi}(t) | \mathcal{G}_t) | \mathcal{F}_t) \\ &= \mathbb{E}_N \left(\lim_{h \downarrow 0} h^{-1} \mathbb{E}_N(\boldsymbol{\xi}(t+h) - \boldsymbol{\xi}(t) | \mathcal{G}_t) | \mathcal{F}_t \right) \\ &= \mathbb{E}_N(\boldsymbol{\alpha}^{(N)}(t) | \mathcal{F}_t), \end{aligned}$$

using a version of the dominated convergence theorem for conditional expectations; see, for example, [10], Section 9.7.

Conversely to the above, a given conditional intensity process yields uniquely the law of a corresponding point process, subject to the condition,

which we shall require, that no two points occur simultaneously. See for example [1]. Here, the limiting measure \mathbb{P} will be specified by the initial distribution ψ and by a conditional intensity process $\boldsymbol{\beta}(t)$ which we will construct; $\boldsymbol{\beta}$ will be adapted to F and represents the limit of the processes $\boldsymbol{\beta}^{(N)}$.

We assume that $\psi_N \rightarrow \psi$ weakly; hence the marginal distribution of $\mathbf{x}(0)$ on \mathbb{Z}_+^K under ψ_N converges weakly to that under ψ . Thus, since we use the discrete topology on \mathbb{Z}_+^K , $\text{Var}_0(\mathbb{P}_N, \mathbb{P}) \rightarrow 0$. (Here Var_0 is the variation distance between the two measures restricted to \mathcal{F}_0 , the σ -algebra containing information only about the initial state $\mathbf{x}(0)$ of the K tagged channels). Theorem 1 of [6] then shows that a sufficient condition for $\text{Var}_t(\mathbb{P}_N, \mathbb{P}) \rightarrow 0$ as $N \rightarrow \infty$ is that

$$(3.4) \quad \int_0^t \mathbb{E}_N |\boldsymbol{\beta}_{k,l}^{(N)}(s) - \beta_{k,l}(s)| ds \rightarrow 0 \quad \text{as } N \rightarrow \infty,$$

for all k, l , where $\boldsymbol{\beta}$ is the conditional intensity of the point process $\boldsymbol{\xi}$ with respect to $(\Omega, \mathbb{P}, \mathcal{F}, F)$.

Since $\boldsymbol{\beta}$ is F -adapted, we have, by (3.3) and the \mathcal{L}^1 -contraction property of conditional expectation (see again [10], Section 9.7), that

$$\begin{aligned} \int_0^t \mathbb{E}_N |\boldsymbol{\beta}_{k,l}^{(N)}(s) - \beta_{k,l}(s)| ds &= \int_0^t \mathbb{E}_N \left| \mathbb{E}_N [\alpha_{k,l}^{(N)}(s) - \beta_{k,l}(s) | \mathcal{F}_s] \right| ds \\ &\leq \int_0^t \mathbb{E}_N |\alpha_{k,l}^{(N)}(s) - \beta_{k,l}(s)| ds. \end{aligned}$$

Hence it suffices to construct a nonnegative F -adapted process $\boldsymbol{\beta}$ such that

$$(3.5) \quad \int_0^t \mathbb{E}_N \left| \mathbb{E}_N [\alpha_{k,l}^{(N)}(s) - \beta_{k,l}(s) | \mathcal{F}_s] \right| ds \rightarrow 0 \quad \text{as } N \rightarrow \infty,$$

for all k, l .

Consider a channel at node i whose queue length is x , while the “environment” is $\underline{\mathbf{u}} \in \mathcal{W}^J$. We define the quantity $c_i(\underline{\mathbf{u}}, x)$ by

$$c_i(\underline{\mathbf{u}}, x) = u_i(x)^{m-1} + u_i(x)^{m-2}u_i(x+1) + \cdots + u_i(x+1)^{m-1}$$

which, provided $u_i(x) > u_i(x+1)$, is equal to

$$\frac{u_i(x)^m - u_i(x+1)^m}{u_i(x) - u_i(x+1)}.$$

Here the numerator is the probability, if the overall state of the network is described by $\underline{\mathbf{u}}$, that out of m randomly chosen queues at node i , the shortest has length x , and the denominator is the proportion of queues at node i with length x . Hence $N^{-1}c_i(\underline{\mathbf{u}}, x)$ may be interpreted as the probability that a new customer arriving at the node chooses this particular channel to join.

Note that for all $x \in \mathbb{Z}_+^K$ and $\underline{\mathbf{u}}, \underline{\mathbf{u}}' \in \mathcal{W}^J$,

$$(3.6) \quad |c_i(\underline{\mathbf{u}}, x)| \leq m$$

and

$$(3.7) \quad |c_i(\underline{\mathbf{u}}, x) - c_i(\underline{\mathbf{u}}', x)| \leq (x+1)m^2 d(\underline{\mathbf{u}}, \underline{\mathbf{u}}').$$

We now consider the conditional intensities of the three types of component process $\xi_{k,l}$, depending on which of the subscripts k and l are zero.

First consider the case $1 \leq k, l \leq K$, so that we are interested in transfers from the k th tagged channel to the l th tagged channel. At time t , the instantaneous rate at which departures destined for node $i(l)$ occur at the k th tagged channel is $I\{x_k(t) > 0\}\mu_{i(k)}p_{i(k)i(l)}$, so that we have

$$\alpha_{k,l}^{(N)}(t) = I\{x_k(t) > 0\}\mu_{i(k)}p_{i(k)i(l)}N^{-1}c_{i(l)} \times (\mathbf{r}(t) - N^{-1}\mathbf{e}_{i(k)}(x_k(t)), x_l(t) - I\{k = l\}).$$

[Here and below we write $\mathbf{e}_i(n)$ for the vector in \mathscr{W}^J whose only nonzero entry is the (i, n) entry, which is 1.] Then

$$(3.8) \quad |\alpha_{k,l}^{(N)}(t)| \leq N^{-1}m \max_i \mu_i \max_{i,j} p_{ij}$$

$$(3.9) \quad \rightarrow 0 \quad \text{as } N \rightarrow \infty,$$

so putting

$$(3.10) \quad \beta_{k,l}(t) = 0$$

gives (3.5) for this case.

Next take $1 \leq k \leq K, l = 0$. We have

$$\alpha_{k,0}^{(N)}(t) = I(x_k(t) > 0)\mu_{i(k)} - \sum_{l=1}^K \alpha_{k,l}^{(N)}(t),$$

so that putting

$$(3.11) \quad \beta_{k,0}(t) = I(x_k(t) > 0)\mu_{i(k)}$$

gives (3.5) again.

The case $k = 0, 1 \leq l \leq K$ is the most difficult; it is when considering arrivals at tagged channels which come from outside the set of tagged channels that the effect of the environment is most greatly felt. We have

$$\alpha_{0,l}^{(N)}(t) = \lambda_{i(l)}c_{i(l)}(\mathbf{r}(t), x_l(t)) + \sum_j \left\{ \left(1 - \frac{\#\{k: i(k)=j\}}{N} \right) \sum_{n=1}^{\infty} \left[r_j(n,t) - r_j(n+1,t) \right] \times \mu_j p_{ji(l)}c_{i(l)}(\mathbf{r}(t) - N^{-1}\mathbf{e}_j(n), x_l(t)) \right\}.$$

Let \mathbf{g} be the point at which the marginal distribution of ψ on \mathscr{W}^J is concentrated. We will set

$$(3.12) \quad \beta_{0,l}(t) = \left[\lambda_{i(l)} + \sum_j u_j(1, t, \mathbf{g})\mu_j p_{ji(l)} \right] c_{i(l)}(\mathbf{u}(t, \mathbf{g}), x_l(t)),$$

where $\mathbf{u}(t, \mathbf{g}) = \{u_j(n, t, \mathbf{g}), 1 \leq j \leq J, n \in \mathbb{Z}_+\}$ is defined by Theorem 2.1(i).

Using (3.6) and (3.7) one can find constants D_1, D_2 and D_3 , depending on $m, \underline{\lambda}, \underline{\mu}$ and \mathbf{P} but not on N , such that

$$\left| \alpha_{0,l}^{(N)}(t) - \beta_{0,l}(t) \right| \leq D_1 \min \left\{ D_2, x_l(t) d(\underline{\mathbf{r}}(t), \underline{\mathbf{u}}(t, \underline{\mathbf{g}})) \right\} + \frac{D_3}{N}.$$

Thus

$$\begin{aligned} & \limsup_{N \rightarrow \infty} \int_0^t \mathbb{E}_N |\alpha_{0,l}^{(N)}(s) - \beta_{0,l}(s)| ds \\ (3.13) \quad & \leq D_1 \limsup_{N \rightarrow \infty} \int_0^t \mathbb{E}_N \min \left(D_2, x_l(s) d(\underline{\mathbf{r}}(s), \underline{\mathbf{u}}(s, \underline{\mathbf{g}})) \right) ds \\ & \leq D_1 \limsup_{N \rightarrow \infty} \int_0^t \left[y \mathbb{E}_N (d(\underline{\mathbf{r}}(s), \underline{\mathbf{u}}(s, \underline{\mathbf{g}}))) + D_2 \mathbb{P}_N (x_l(s) > y) \right] ds \end{aligned}$$

for any y . The arrival process at channel l may be dominated for all N by a Poisson process of rate $m(\lambda_{i(l)} + \sum_j \mu_j p_{ji(l)})$, and the distribution of $x_l(0)$ under ψ_N converges to that under ψ , so we have

$$\mathbb{P}_N (x_l(s) > y) \rightarrow 0 \quad \text{as } y \rightarrow 0,$$

uniformly in N and in $s \in [0, t]$. Also, from Theorem 2.1(ii) and the definition (2.2) of the metric d ,

$$\mathbb{E}_N d(\underline{\mathbf{r}}(s), \underline{\mathbf{u}}(s, \underline{\mathbf{g}})) \rightarrow 0 \quad \text{as } N \rightarrow \infty,$$

uniformly in $s \in [0, t]$. This shows that the RHS of (3.13) is 0, so that (3.5) holds in this final case also.

Thus, since β defined by (3.10), (3.11) and (3.12) is nonnegative and F -adapted and satisfies (3.5) for all k and l , the first part of the theorem is proved.

For the second part, note that, under \mathbb{P} , the instantaneous transition rates at time t of the process ξ , and hence also of the process \mathbf{x} , can be read off immediately from the vector $\beta(t)$. However, $\beta(t)$ is a function of $\mathbf{x}(t)$ —in particular it depends on $\{\mathbf{x}(s), 0 \leq s \leq t\}$ only through $\mathbf{x}(t)$ —so \mathbf{x} is Markov.

Further, note that $\beta_{k,l}(t)$ is identically zero if $k > 0$ and $l > 0$ and that $\beta_{k,0}(t)$ and $\beta_{0,k}(t)$ depend only on $x_k(t)$. Conversely, $x_k(t)$ is a function only of $\xi_{k,l}(t)$ and $\xi_{l,k}(t)$, $0 \leq l \leq K$. So in fact the instantaneous transition rates of x_k at time t depend on $\{\mathbf{x}(s), 0 \leq s \leq t\}$ only through $x_k(t)$, and so if the components of \mathbf{x} are independent at time 0, then their paths thereafter are also independent. \square

4. Systems in equilibrium. If the networks considered in the previous section are each started in equilibrium, we can describe more precisely the limiting point process arising from Theorem 3.1.

We will require that no node is *overloaded*. Assume that the matrix $\mathbf{I} - \mathbf{P}$ is invertible. (A physical interpretation of this condition is that, almost surely, every task entering the network eventually leaves the network). Define the vector $\underline{\rho} = (\rho_1, \dots, \rho_J)$ by

$$\underline{\rho} = \underline{\lambda}(\mathbf{I} - \mathbf{P})^{-1}.$$

Then we have, for all i ,

$$(4.1) \quad \rho_i = \lambda_i + \sum_{1 \leq j \leq J} \rho_j P_{ji},$$

and if

$$(4.2) \quad \underline{\rho} < \underline{\mu} \text{ (i.e., } \rho_j < \mu_j \text{ for all } j),$$

then the networks are said to be nonoverloaded, and $N\rho_i$ may be interpreted as the “effective arrival rate” in equilibrium at node j in the N th network, including arrivals from inside as well as outside the network.

In [8] the following result is established concerning the equilibrium behavior of the networks considered.

THEOREM 4.1. *If (4.2) holds, then:*

(i) *There exists a unique fixed point $\underline{\mathbf{a}}$ in \mathcal{U}^J of the system (2.3)–(2.6), that is, such that $\underline{\mathbf{u}}(t, \underline{\mathbf{a}}) = \underline{\mathbf{a}}$ for all t , and $\underline{\mathbf{a}}$ is given by*

$$(4.3) \quad a_j(n) = \left(\frac{\rho_j}{\mu_j} \right)^{(m^n - 1)/(m - 1)}.$$

(ii) *The Markov process $\underline{\mathbf{r}}_N(t)$ is positive recurrent for all N , and so has a unique invariant probability distribution π_N for each N .*

(iii) *$\pi_N \rightarrow \delta_{\underline{\mathbf{a}}}$ weakly as $N \rightarrow \infty$, where $\delta_{\underline{\mathbf{a}}}$ is the probability measure concentrated at the fixed point $\underline{\mathbf{a}}$.*

So from here on we assume that (4.2) holds and discuss the situation in which, for each N , the distribution under ψ_N of the initial state $\mathbf{r}(0)$ on \mathcal{U}_N^J is the equilibrium distribution π_N given by Theorem 4.1(ii). Then the sequence ψ_N has a weak limit ψ under which $\mathbf{r}(0)$ is equal to $\underline{\mathbf{a}}$ with probability 1, and we can apply Theorem 3.1. Since $\underline{\mathbf{u}}(\underline{\mathbf{a}}, t) = \underline{\mathbf{a}}$ for all t , (3.12) becomes

$$(4.4) \quad \beta_{0,l}(t) = \rho_{i(l)} c_{i(l)}(\underline{\mathbf{a}}, x_l(t)),$$

using (4.1). For $m = 1$, $c_{i(l)}(\cdot) = 1$ identically, and for $m > 1$ we have

$$c_i(\underline{\mathbf{a}}, x) = \frac{(\rho_i/\mu_i)^{m(m^x-1)/m-1} (1 - (\rho_i/\mu_i)^{m^{x+1}})}{(\rho_i/\mu_i)^{m^x-1/m-1} (1 - (\rho_i/\mu_i)^{m^x}}.$$

As before, $\beta_{k,l}(t) = 0$ for $k, l \geq 1$, and $\beta_{k,0}(t) = I\{x_k(t) > 0\} \mu_{i(k)}$ for $k \geq 1$. Thus under the limiting measure \mathbb{P} , the queue length processes at the tagged channels are independent time-homogeneous Markov birth-and-death processes. From the fixed point result in Theorem 4.1(i), it follows that the process at a channel belonging to node i has an equilibrium distribution under which the probability of the queue length being n or longer is $\{a_i(n), n \in \mathbb{Z}_+\}$.

In the case $m = 1$, the queue behaves simply as an $M/M/1$ queue, whose arrival rate does not depend on the queue length. For $m > 1$, the arrival rate decreases as the length increases.

The arrival rate when the queue is empty is

$$\frac{1 - (\rho_i/\mu_i)^m}{1 - (\rho_i/\mu_i)},$$

which increases as m increases. Thus the average length of an idle period of the queue decreases as m increases; the overall intensity of arrivals and the average service time stay constant, so the average length of a busy period decreases also.

We now discuss specifically the point process of arrivals at such a queue. (If the queue is started in equilibrium, then by a reversibility argument one can show that this has the same distribution as the point process of departures). If $m = 1$, this is simply a Poisson process of rate ρ_i . As m increases, the average intensity remains the same, but the points tend to become more “evenly spread”; the probability of a given time interval containing no points at all decreases, for example. In the limit $m = \infty$, the points of the process are those of a renewal process, whose renewal intervals are the independent sum of a service time which is exponentially distributed with rate μ_i and an arrival time which is exponentially distributed with rate $(1 - \rho_i/\mu_i)^{-1}$. This corresponds to a situation where arrivals at a node choose freely between all channels, and so only ever join empty queues.

This change is illustrated by Figure 1, which shows the results of simulations of sequences of 500 interarrival times in the cases $m = 1$, $m = 2$ and $m = \infty$ for $\rho_i = 0.5$ and $\mu_i = 1$. As m increases, very long interarrival distances become less frequent and tend to bunch together less; the same applies to very short interarrival distances.

A different phenomenon occurs if ρ_i is nearly 1. Then the difference between the Poisson process of arrivals for $m = 1$ and the renewal process of arrivals for $m = \infty$ is very slight and, as noted above, in each case the departure process is distributed exactly as the arrival process, so again is very similar for $m = 1$ and $m = \infty$. However, the process of queue lengths at the channel, which is determined by the *joint* distribution of arrivals and departures, is extremely different: for example, for $m = 1$ its stationary distribution is geometric with mean $\rho_i/(1 - \rho_i)$, while for large m the queue length hardly ever exceeds 1.

We can alternatively consider the network from the viewpoint of a particular task progressing through it. Since the routing is Markovian, its route may as well be considered fixed as soon as it enters the network: the route may be taken to include both the order in which the task visits nodes and the particular channels it inspects on each arrival at a node.

The above observation that the queue length processes at the tagged channels are independent in the limit can then be interpreted as follows: as $N \rightarrow \infty$, we approach a situation in which the waiting times of the task at each stage of the route are independent, and where the queue lengths of the m channels that the task inspects at node i are independently drawn from the distribution represented by $a_i(n)$, $n \in \mathbb{Z}_+$. Thus the probability that the task joins a queue of length n at node i is $a_i(n)^m - a_i(n+1)^m$, which, from (4.3),

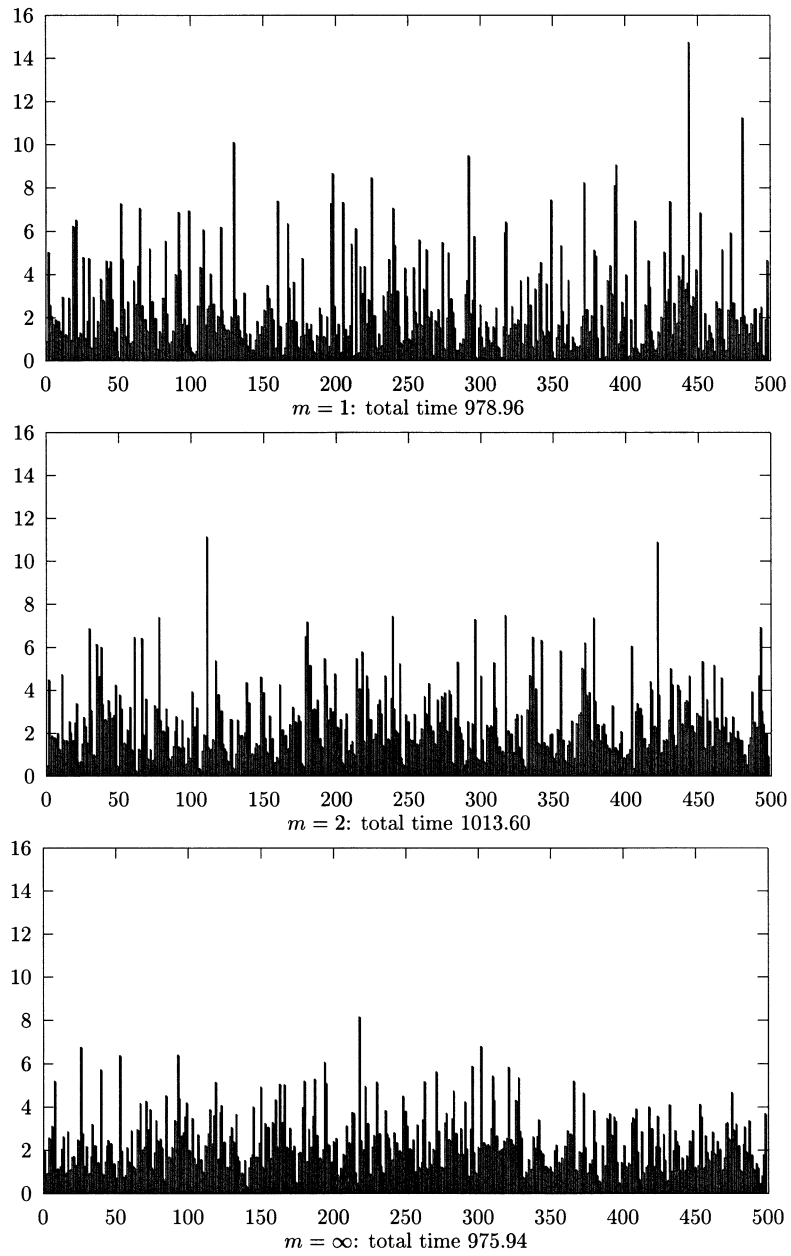


FIG. 1. Simulations of 500 successive interarrival times for $\rho = 0.5$, $\mu = 1$, various m .

can be written as

$$\left(\frac{\rho_i}{\mu_i}\right)^{m(m^n-1)/(m-1)} \left(1 - \left(\frac{\rho_i}{\mu_i}\right)^{m^{n+1}}\right),$$

decaying superexponentially in n for $m > 1$.

5. Extensions. An analogous result to Theorem 3.1 holds if we remove the condition that the marginal distribution of $\mathbf{r}(0)$ on \mathcal{Z}^J under the limiting initial condition ψ is concentrated at a single point. Then, however, we must replace (3.12) by

$$\alpha_{0,l}(t) = \left[\lambda_{i(l)} + \sum_j u_j(\mathbf{1}, t, \mathbf{r}(0)) \mu_j p_{ji(l)} \right] c_{i(l)}(\mathbf{u}(t, \mathbf{r}(0)), x_l(t))$$

and

$$\beta_{0,l}(t) = \mathbb{E}[\alpha_{0,l}(t) | \mathcal{F}_t].$$

Now part (ii) of Theorem 3.1 fails, since the observed paths of $\{\mathbf{x}(s), 0 \leq s \leq t\}$ provide information about the “hidden” initial environment $\mathbf{r}(0)$ and hence about $\mathbf{r}(t)$; so in general \mathbf{x} is no longer Markov and the components of $\mathbf{x}(t)$ are not independent even if those of $\mathbf{x}(0)$ are. One can formulate an interesting filtering problem concerning the estimation of the initial environment $\mathbf{r}(0)$ given the observed paths $\{\mathbf{x}(s), 0 \leq s \leq t\}$ of the queue length processes at the tagged channels.

Following more closely the approach of Brown and Pollett in [2], we can consider letting J , the number of nodes, rather than N , the number of channels at each node, tend to infinity. In [2], this limit is considered for single-class Markovian queueing networks with state-dependent service rates, and, under various conditions, bounds are derived for the variation distance between the equilibrium arrival process at a node and a Poisson process, which tend to 0 as $J \rightarrow \infty$. For example, an appropriate condition for standard Jackson networks is that

$$\sum_{j=1}^J (\mu_j p_{ji})^2 \rightarrow 0 \quad \text{as } J \rightarrow \infty.$$

It seems likely that similar results hold also in the situations we have considered above. However, the networks considered in [2] all have the property that the equilibrium distribution of the state of the network has a product form, and this is an important element of the methods used there. In the networks we have considered here, the lengths of the queues at channels at the same or at different nodes are not in general independent for finite N , even in equilibrium, and it seems that different methods will be needed to establish such results in this case.

Acknowledgments. I thank a referee for detailed comments which led to considerable improvements in the rigor and presentation of the paper and Dr. Yuri Suhov for many valuable conversations during the course of this work.

REFERENCES

- [1] AALEN, O. (1978). Nonparametric inference for a family of counting processes. *Ann. Statist.* **6** 701–726.
- [2] BROWN, T. C. and POLLETT, P. K. (1982). Some distributional approximations in Markovian queueing networks. *Adv. Appl. Probab.* **14** 654–671.
- [3] DALEY, D. J. and VERE-JONES, D. (1988). *An Introduction to the Theory of Point Processes*. Springer, New York.
- [4] JACKSON, J. R. (1957). Networks of waiting times. *Oper. Res.* **5** 518–527.
- [5] JACKSON, J. R. (1965). Jobshop-like queueing systems. *Management Sci.* **10** 131–142.
- [6] KABANOV, Y. M. and LIPTSER, R. S. (1983). Convergence of the distributions of multivariate point processes. *Z. Wahrsch. Verw. Gebiete* **63** 475–485.
- [7] KELLY, F. P. (1979). *Reversibility and Stochastic Networks*. Wiley, New York.
- [8] MARTIN, J. B. and SUHOV, Y. M. (1999). Fast Jackson networks. *Ann. Appl. Probab.* **9** 854–870.
- [9] VVEDENSKAYA, N. D., DOBRUSHIN, R. L. and KARPELEVICH, F. I. (1996). Queueing system with selection of the shortest of two queues: an asymptotic approach. *Problems Inform. Transmission* **32** 15–27.
- [10] WILLIAMS, D. (1991). *Probability with Martingales*. Cambridge Univ. Press.

DÉPARTEMENT DE MATHÉMATIQUES
ET D'INFORMATIQUE
ECOLE NORMALE SUPÉRIEURE
45 RUE D'ULM
75230 PARIS CÉDEX 05
FRANCE
E-MAIL: james.martin@ens.fr