

Diffusion Approximations for a Single Node Accessed by Congestion-Controlled Sources

Atanu Das and R. Srikant

Abstract—We consider simple models of congestion control in high-speed networks, and develop diffusion approximations which could be useful for resource allocation. We first show that, if the sources are ON-OFF type with exponential ON and OFF times, then, under a certain scaling, the steady-state distribution of the number of active sources can be described by a combination of two appropriately truncated and renormalized normal distributions. For the case where the source arrival process is Poisson and the service times are exponential, the steady-state distribution consists of appropriately normalized and truncated Gaussian and exponential distributions. We then consider the case where the arrival process is a general renewal process with finite coefficient of variation and service-time distributions that are phase type, and show the impact of these distributions on the steady-state distribution of the number of sources in the system. We also establish an insensitivity to service-time distribution when the arrival process is Poisson. We use these results to relate the capacity of a bottleneck node to performance measures of interest for best effort traffic, such as the mean file transfer time or probability of congestion.

I. INTRODUCTION

TRADITIONALLY telephone networks have provided a guaranteed bandwidth to each call, and when the capacity of a link is exceeded, further calls would be rejected. When the arrival process is Poisson, traffic engineering for such networks can be performed using the well-known Erlang-B formula [6]. The Erlang-B formula is extremely useful due to the well-known insensitivity to the service-time distribution. Often, in telephone networks, routing phenomena such as overflow do not preserve the Poisson nature of arrivals at a particular link. Under these circumstances, to obtain blocking probabilities or other performance measures, one has to take into account the statistics of the arrival process and service times. A precise computation is extremely difficult. An appealing alternative is to use heuristics obtained from weak convergence theory and diffusion limits to obtain approximations to the performance measures of interest [40]. In fact, diffusion limits can also be useful in predicting the behavior of the system over finite-time intervals (not just steady state) [32], [33].

Even for the simple $M/M/s/0$ loss models, diffusion approximations are useful in providing simple rules-of-thumb re-

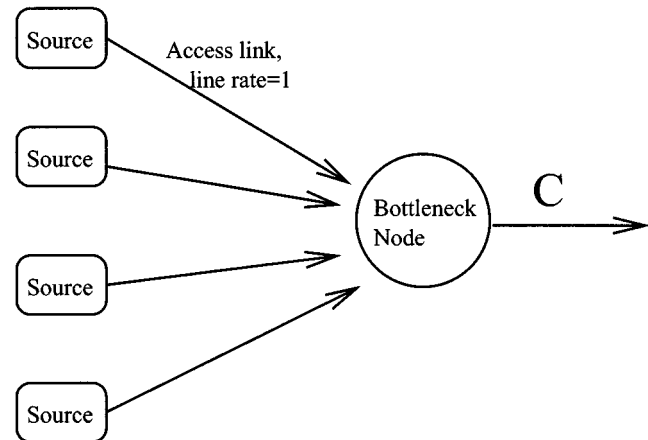


Fig. 1. Instantaneous model of congestion-controlled best effort traffic.

lationships among between offered load, capacity, and performance measures such as blocking probability. For example, the results in [40] and [32] indicate that, under critical loading (i.e., the offered load is roughly on the order of the capacity), the exact values of the offered load and capacity are not important. For the purposes of traffic engineering, the factor that determines the blocking probability is the difference between the offered load and capacity, measured in units of the square root of the offered load (or capacity). Thus, if the offered load is increased, using the diffusion approximation, one can approximately compute the required increase in provisioned capacity to maintain the same blocking probability.

In this paper, we develop similar diffusion approximations for congestion-controlled systems. The basic model is illustrated in Fig. 1. We consider a single node with capacity C (measured possibly in bits per second). At any instant in time, the node is accessed by several sources performing such tasks as Web browsing, file transfer, etc. We assume that the access lines for each sources have a capacity of one unit each. Thus, when the total number of sources using the node is less than or equal to C , the sources are *access limited*, i.e., the response of the network is limited by the access link capacities. However, when the number of active sources is greater than C , then we assume that there is some congestion-control mechanism that attempts to divide the capacity fairly among the competing sources. This model was proposed and analyzed for Web traffic in [16].

The interest in such a model is due to the dramatic development of best effort services over the last few years. Unlike traditional telephone networks, best effort traffic is not allocated a fixed bandwidth in the network. Rather, the available capacity is divided among various competing sources in a fair manner. The

Manuscript received July 22, 1998; revised May 10, 1999 and December 10, 1999. Recommended by Associate Editor, L. Dai. This work was supported by NSF CAREER Award NCR 9701525, NSF Grant ANIR 981370, and a grant from Nokia Research Center.

A. Das is with the Coordinated Science Laboratory and the Department of Electrical and Computer Engineering, University of Illinois, Urbana, IL 61801 USA (e-mail: adas@uiuc.edu).

R. Srikant is with the Coordinated Science Laboratory and the Department of General Engineering, University of Illinois, Urbana, IL 61801 USA (e-mail: rsrikant@uiuc.edu).

Publisher Item Identifier S 0018-9286(00)09448-4.

two most common approaches to this are TCP and ATM ABR [38]. In the Internet, the well-known TCP protocol (see [18] for the original congestion avoidance scheme in TCP and [23], and references within, for subsequent modifications) is a feedback-control mechanism that probes the network for available capacity, and increases or decreases the rate at which sources are sending traffic into the network, depending upon the congestion in the network. In the context of asynchronous transfer mode (ATM) networks, the available bit-rate (ABR) service performs the congestion-control function, again through a feedback control mechanism [4], [30], [19], [1], [2], [5]. In either case, typically, the sharing mechanism that allocates the node's capacity to the competing sources is not efficient. For instance, [16] presents an analytical model to quantify this inefficiency in the case of TCP. According to their model, when the number of active sources is greater than C , a fraction (denoted by β) of the total capacity C is shared equally by the active sources. Thus, the factor β will depend on the feedback control mechanism that is used to perform the congestion control. We assume this model in this paper. The intuitive reason for such a model is as follows: consider a single TCP source which increases its window size until it detects congestion due to lost packets, then drops its window size, and repeats this cycle. Thus, instead of operating at the node's data rate, the rate at which a TCP source generates data follows a cycle starting from near zero, and reaching a peak rate which is larger than the rate that can be sustained by its bottleneck node. Thus, the *goodput*, i.e., the rate of successful data transfer, seen by the source would be smaller than the bottleneck node's capacity. We refer the interested reader to [16] where extensive simulations have been used to justify this model.

From this congestion-control scenario, we derive two mathematical models for the congestion-controlled best effort traffic. In the first one, we assume that the number of sources accessing the network is fixed, and that the sources switch from an active (ON) state to an inactive (OFF) state and vice versa according to a Markov chain. This model is appropriate in a scenario where the traffic is predominantly Web browsing, and was considered in [16]. We develop a diffusion approximation for this model by scaling the bottleneck capacity and the number of sources appropriately. The steady-state distribution of the number of ON sources in this model is known to be independent of the ON and OFF time distributions [16]. Thus, our diffusion approximation is primarily useful in understanding the relationship among the capacity, offered load, efficiency factor β , and performance measures, much like the diffusion approximation in the case of the classical Erlang $M/M/s/0$ loss model.

In the second model, we assume that the number of sources in the system is not a constant. We suppose that the sources arrive at the system according to some stationary renewal process with finite interarrival time variance. We believe that this model is more generally applicable in situations where the traffic consists of a mixture of Web browsing, file transfers, and other possible uses for the Internet. Each source brings a certain amount of work to the network. This work can be thought of as the file size for a file transfer. We assume that file sizes are distributed independently according to a phase-type distribution, which we describe precisely later.

The above models are not intended to capture the complex dynamics of congestion-control schemes such as TCP. The idea is to simply use the fact that any congestion control attempts to fairly distribute the available bandwidth, even though the actual implementation may fall short of this goal. Our aim is to build a model that can be useful for network provisioning, not to precisely model the dynamics of the congestion-control mechanism. It is also interesting to note that, in the case of ON-OFF models, the ON time is a function of the congestion in the network. This is an appropriate model for Web traffic only. The rationale for this is as follows [16]: a user first downloads a Web page, and waits for the page to be fully downloaded. After this download is completed, the user may pause for a further amount of time, and follow a hyperlink on the downloaded Web page, which results in a new data transfer. This model may not be applicable more generally, and therefore, we consider the second model.

For both models, we consider two performance measures. The first is the mean file transfer time. When the capacity of the bottleneck node is infinity, the mean file transfer time would simply be the mean file size divided by the speed of the access link. Thus, a relevant measure of the quality of service (QoS) perceived by a user is the mean file transfer time (or the increase in it) due to the finite capacity of the bottleneck link. The second performance measure that we consider is the probability of congestion, i.e., the probability that the sum of the access link speeds of the active sources exceeds the capacity of the bottleneck link. As mentioned earlier, for the ON-OFF source model, both of these performance measures are independent of the ON and OFF time distributions. However, in the case of the renewal arrival process model, these performance measures vary with the interarrival and service-time (file-size) distributions. In general, this relationship can be complex and difficult to characterize exactly. Thus, we provide simple approximations of these performance measures that relate the values of performance measures under general distributions to the values of the performances measures for the case of Poisson arrivals and exponential service times. The idea here is that the solution for the case of Poisson arrivals and exponential service times is easy to obtain by solving a simple birth-death Markov chain, and thus, the simple approximation can be exploited to calculate these performance measures for general distributions. Once we have such a simple relationship between the traffic parameters and the performance measures, traffic engineering is straightforward. One can simply increase or decrease the capacity of the bottleneck node until the desired performance is attained.

The rest of the paper is organized as follows. In Section II, we consider both the ON-OFF and the $M/M/C$ models, and derive the diffusion approximation for both cases. Section III derives the diffusion approximation for the $GI/M/C$ model, and obtains explicit expressions for the stationary distribution. In Section IV, we prove the insensitivity of the stationary distribution to the service-time distribution in the $M/GI/C$ model. Section V deals with the $GI/GI/C$ model, and lays the groundwork for the weak convergence proof in Appendix III by computing the drift vector and diffusion coefficient matrix for the limiting diffusion process. In Section VI, we use the diffusion limits to obtain expressions relating the traffic statistics to the

QoS delivered to the user by the network. Section VII presents numerical results to support the diffusion approximations, and concluding remarks are provided in Section VIII. Appendixes I and II provide alternate derivations of the stationary distributions in the case of ON-OFF models and $M/M/C$ models, using Stirling's formula and the central limit theorem.

II. EXPONENTIAL INTERARRIVAL AND SERVICE TIMES

A. ON-OFF Sources

Consider S ON-OFF sources which switch between an active (ON) state and an inactive (OFF) state. Let us suppose that the sources access a bottleneck node with transmission rate C bits/s, and that each source accesses this node via an access link whose transmission rate is 1 bit/s. In other words, when the number of ON sources in the system is less than or equal to C , each source is served at rate 1. Let $N(t)$ be the number of sources that are in the ON state at time t . We assume that a modified form of processor sharing is in effect, i.e., when the number of sources exceeds C , there is some form of congestion-control mechanism such as TCP [16] or ATM ABR rate control [2] that divides the available capacity C to all of the sources equally. However, most congestion-control schemes are not perfect, and therefore, as will be made precise later, we introduce an efficiency factor to account for this. This model was used by Heyman *et al.* [16], and we will call this type of service discipline *modified processor sharing* (MPS). We assume that the OFF times are exponential with mean $1/\lambda$. Once a source enters the ON state, it remains in this state until it receives an exponentially distributed amount of service whose mean is $1/\mu$ bits. Since the access rate is 1, $1/\mu$ would be the mean time in the ON state if the number of sources is less than C . Thus, we will refer to $1/\mu$ as the *nominal* mean ON time. In general, the time spent in the ON state by a source will also depend on the number of other sources that are ON.

We emphasize that this model assumes a fixed number of sources accessing the node, with some sources being ON and the others being OFF at any time. In a Web-browsing application, a user typically downloads a Web page, and looks at it for a while before generating another request for a new download. Thus, since we have assumed a fixed number of total sources, when there are many sources in the ON state, the rate at which new files are generated also slows down. We refer the reader to [16], [3] for a justification of this model using actual traffic traces.

Let q_{ij} be the transition rate from state $\{N(t) = i\}$ to $\{N(t) = j\}$ in the associated birth-death process. Then, q_{ij} is given by

$$q_{i,i+1} = (S - i)\lambda, \quad i < S \quad (1)$$

and

$$q_{i,i-1} = \begin{cases} i\mu, & \text{if } i \leq C \\ \beta C\mu, & \text{if } i > C \end{cases} \quad (2)$$

where $\beta < 1$ is a factor denoting the efficiency of the congestion-control scheme. Note that, with $C = 1$ and $\beta = 1$, this becomes the well-known processor-sharing service discipline which has been studied extensively; see [24, ch. 4.4] and references within.

Let us consider the following diffusion scaling for the above process:

$$S_n = n, \quad C_n = n\rho + \sqrt{n\rho}\gamma, \quad \beta_n = 1 - \frac{\theta}{\sqrt{n\rho}} \quad (3)$$

and define

$$x_n(t) = \frac{N(t) - n\rho}{\sqrt{n}}$$

where $\rho \triangleq \lambda/(\lambda + \mu)$ is the offered load per source. We are interested in the behavior of $x_n(t)$ as $n \rightarrow \infty$. The scaling for S_n and C_n reflects the fact that we are interested in studying the behavior of systems where the number of sources is large and the capacity of the system is increased such that $S_n\rho/C_n \rightarrow 1$. Such a scaling is referred to as the *critical-loading* regime. For systems with a large number of sources, other traffic regimes of interest are the *heavy-loading* regime where $\lim_{n \rightarrow \infty} (S_n\rho/C_n) > 1$, and the *light-loading* regime where $\lim_{n \rightarrow \infty} (S_n\rho/C_n) < 1$. Service providers typically operate their network in the *critical-loading* regime, and thus, it is of most interest for traffic engineering purposes. We only study the *critical-loading* regime here. The reason for the particular choice of the scaling β_n is addressed later in this section in Remark 1. With $\beta_n = 1$, the above scaling is well known for infinite-server [17], [9], [39], [13], [14] and loss models [10], [40], [31], [32].

Let us first consider the case $N(t) < C$, which is equivalent to $x(t) < \gamma\sqrt{\rho}$. Then, the drift and diffusion coefficients are given by

$$\begin{aligned} m(x) &\triangleq \lim_{\delta \rightarrow 0} E\left(\frac{x_n(t+\delta) - x_n(t)}{\delta} \middle| x_n(t) = x\right) \\ &= -(\lambda + \mu)x \end{aligned} \quad (4)$$

and

$$\begin{aligned} \sigma^2(x) &\triangleq \lim_{\delta \rightarrow 0} E\left(\frac{(x_n(t+\delta) - x_n(t))^2}{\delta} \middle| x_n(t) = x\right) \\ &= 2\rho\mu \end{aligned} \quad (5)$$

respectively. For the case $N(t) > C$ or, equivalently $x(t) > \gamma$, the drift is given by

$$m(x) = -\lambda x + \mu\sqrt{\rho}(\theta - \gamma) \quad (6)$$

and the diffusion coefficient is the same as before.

The steady-state density for this process is obtained by solving (see, for example, [21])

$$\frac{1}{2} \frac{d^2}{dx^2} (\sigma^2(x)p(x)) = \frac{d}{dx} (m(x)p(x)). \quad (7)$$

For $x < \gamma\sqrt{\rho}$, solving the above equation as in [21] yields

$$\rho\mu p(x)e^{((\lambda+\mu)x^2/2\rho\mu)} = C_1 \int_0^x e^{((\lambda+\mu)y^2/2\rho\mu)} dy + C_2$$

for some constants C_1 and C_2 . Since $p(x)$ should be integrable about $x = -\infty$, $C_1 = 0$. Thus,

$$p(x) = K_1 e^{-((\lambda+\mu)x^2/2\rho\mu)}, \quad \text{for } x < \gamma\sqrt{\rho} \quad (8)$$

for some constant K_1 . For $x > \gamma\sqrt{\rho}$, solving (7) gives

$$\begin{aligned} & \mu p(x) e^{(1/\rho\mu)((\lambda x^2/2) + \mu\sqrt{\rho}\lambda x - \mu\theta x)} \\ &= C_3 \int_{2\gamma}^x e^{(1/\rho\mu)((\lambda y^2/2) + \mu\sqrt{\rho}\lambda y - \mu\sqrt{\rho}\theta y)} dy + C_4. \end{aligned}$$

Again, since $p(x)$ should be integrable about $x = -\infty$, we obtain

$$p(x) = K_2 e^{(-1/\rho\mu)((\lambda x^2/2) + \mu\sqrt{\rho}\lambda x - \mu\sqrt{\rho}\theta x)}. \quad (9)$$

To explicitly show that $p(x)$ has the form of two different normal densities in the regions $x < \gamma\sqrt{\rho}$ and $x > \gamma\sqrt{\rho}$, we rewrite (8) and (9) as

$$p(x) = K_1 e^{-(x^2/2\rho(1-\rho))}, \quad \text{for } x < \gamma\sqrt{\rho} \quad (10)$$

$$p(x) = K_2 e^{-(1/2(1-\rho))(x - ((\theta - \gamma)(1 - \rho)/\sqrt{\rho}))^2}, \quad (11)$$

$$\text{for } x > \gamma\sqrt{\rho}$$

where, by abusing notation, we use K_2 to denote a new constant. Then, for $x < \gamma\sqrt{\rho}$, the normal density is centered at zero, while for $x > \gamma\sqrt{\rho}$, the normal density is centered at $((\theta - \gamma)(1 - \rho)/\sqrt{\rho})$.

To solve for K_1 and K_2 , we have the following two conditions: the continuity condition

$$\begin{aligned} & K_1 e^{-((\gamma\sqrt{\rho})^2/2\rho(1-\rho))} \\ &= K_2 e^{-(1/2(1-\rho))(\gamma\sqrt{\rho} - ((\theta - \gamma)(1 - \rho)/\sqrt{\rho}))^2} \end{aligned} \quad (12)$$

and the normalization condition

$$\begin{aligned} & K_1 \sqrt{2\pi\rho(1-\rho)} \left(1 - Q\left(\frac{\gamma\sqrt{\rho}}{\sqrt{\rho(1-\rho)}}\right) \right) \\ &+ K_2 \sqrt{2\pi(1-\rho)} Q\left(\frac{\gamma\sqrt{\rho} - \frac{(\theta - \gamma)(1 - \rho)}{\sqrt{\rho}}}{\sqrt{1-\rho}}\right) = 1 \end{aligned} \quad (13)$$

where $Q(x)$ is the complementary cumulative distribution function (ccdf) of the Gaussian random variable $N(0, 1)$. The above discussions lead to the following theorem, whose proof can be found in Appendix III.

Theorem 1: Suppose that $(N(0) - n\rho/\sqrt{n}) \Rightarrow x(0)$; then

$$\frac{N(n\cdot) - n\rho}{\sqrt{n}} \Rightarrow x(\cdot)$$

where \Rightarrow denotes weak convergence in $D[0, \infty)$, the space of all right continuous functions with left limits. The convergence is with respect to the Skorohod's J_1 topology [7]. The limiting process $x(t)$ is described the stochastic differential equation

$$dx = m(x)dt + \sigma dw \quad (14)$$

with $x(0)$ given, $m(x)$ is given by (4) and (6), σ is given by (5), and $w(t)$ is a standard Wiener process. The steady-state distribution of $x(t)$ is given by (10) and (11). \diamond

The above results have a technical gap that we have not addressed. We can show the weak convergence of the process, and we know the stationary distribution of the limiting process. However, this does not immediately mean that the stationary distribution of the n th system converges in the limit. By considering $P(N > C)$, we show that this convergence indeed occurs in Appendix I, by directly working with the Markov chain of the n th system.

B. Poisson Source

Consider a Poisson arrival process with rate λ and exponentially distributed service times with mean $1/\mu$. As before, let the maximum number of sources that can be served by the system without processor sharing taking effect be C . In this case, the transition rates for the birth-death Markov chain describing the number of active sources is given by

$$q_{i,i+1} = \lambda \quad (15)$$

and

$$q_{i,i-1} = \begin{cases} i\mu, & \text{if } i \leq C \\ \beta C\mu, & \text{if } i > C \end{cases} \quad (16)$$

where $\beta < 1$ is the efficiency factor defined earlier.

In this case, we consider the following scaling:

$$\lambda = n\mu, \quad C = n + \gamma\sqrt{n}, \quad \text{and} \quad \beta = 1 - \frac{\theta}{\sqrt{n}}. \quad (17)$$

As before, we define $x(t) = N(t) - n/\sqrt{n}$, where $N(t)$ is the number of active sources at time t . Then, $m(x)$ is given by

$$m(x) = \begin{cases} -\mu x, & \text{if } x < \gamma \\ (\theta - \gamma)\mu, & \text{if } x > \gamma \end{cases} \quad (18)$$

and $\sigma^2(x)$ given by

$$\sigma^2(x) = 2\mu, \quad \forall x. \quad (19)$$

For stability, we require that $\beta_n C_n \mu > n$. This implies that $\gamma - \theta > 0$. Proceeding as before, we get

$$p(x) = \begin{cases} K_1 e^{-x^2/2}, & \text{if } x < \gamma \\ K_2 e^{-(\gamma - \theta)x}, & \text{if } x > \gamma. \end{cases} \quad (20)$$

Thus, in the case of a Poisson source model, the diffusion approximation for the steady-state density function consists of a renormalized normal density and an exponential density.

In this case, using the continuity and normalization conditions as in the previous subsection, we can obtain

$$K_1 = \frac{1}{\sqrt{2\pi(1 - Q(\gamma))} + \frac{e^{-\gamma^2/2}}{\gamma - \theta}} \quad (21)$$

and

$$K_2 = \frac{e^{\gamma^2/2} e^{-\theta\gamma}}{\sqrt{2\pi(1 - Q(\gamma))} + \frac{e^{-\gamma^2/2}}{\gamma - \theta}}. \quad (22)$$

From these expressions, we can compute the steady-state probability that the system is in a congested state, i.e., the number of sources in the system is greater than C , and is given by

$$\begin{aligned} & \lim_{t \rightarrow \infty} P(N(t) > C) \\ & \approx P(x > \gamma) \\ & = \frac{e^{-\gamma^2/2}}{\sqrt{2\pi}(\gamma - \theta)(1 - Q(\gamma)) + e^{-\gamma^2/2}}. \end{aligned} \quad (23)$$

The above derivation can be summarized in the following theorem, the proof of which is a special case of the proof for the $GI/M/C$ model in Appendix III.

Theorem 2: Suppose that $(N(0) - n/\sqrt{n}) \Rightarrow x(0)$; then

$$\frac{N(n \cdot) - n}{\sqrt{n}} \Rightarrow x(\cdot)$$

where $x(t)$ is given the stochastic differential equation (14), $m(x)$ is given by (18), and σ is given by (19). The steady-state distribution of $x(t)$ is given by (20). \diamond

Remark 1: We note that the case $\beta_n \equiv 1$ (or, equivalently, $\theta = 0$) was considered by [15]. In that case, the drift $m(x)$ is continuous. However, when $\theta \neq 0$, $m(x)$ is discontinuous, and we cannot use Stone's theorem [34], [17, Theorem 3.2] to prove the weak convergence of $x_n(t)$ to $x(t)$. Further, in the case $\beta \neq 1$, β_n should be scaled as we have done; otherwise, one cannot establish the desired weak convergence, nor can one directly work with the Markov chain to obtain the stationary distribution as $n \rightarrow \infty$. In fact, it is straightforward to check that the required condition is that the limit $\lim_{n \rightarrow \infty} \sqrt{n}(1 - \beta_n)$ exists. \diamond

The fact that the steady-state distribution $P(N_n(\infty) > C_n)$ of the n th system converges to the corresponding steady-state distribution of $x(t)$ is verified in Appendix II. However, this alone may not be sufficient in practice [28], [26], [27]. A desirable result would be the following: given $\epsilon > 0$, there exists an x such that $P(|x_n(t)| > x) < \epsilon$ for all n, t . For our problem, this is easy to verify as follows. As in [15], for each n , $x_n(t)$ can be stochastically upper bounded by another process obtained by placing an impenetrable lower barrier at $N_n = C_n$. Further, $x_n(t)$ can be stochastically lower bounded by an $M/M/\infty$ system. Since the upper bound is an $M/M/1$ queue and the lower bound is an $M/M/\infty$ system, they are easy to analyze. Starting from an empty system, it is well known that the number of customers in an $M/M/1$ queue stochastically increases to its steady-state distribution. For the $M/M/\infty$ system, starting from an empty system, the number in the system at any time is a Poisson random variable. Using these facts, the desired result can be proved easily.

III. $GI/M/C$ QUEUE WITH MPS

Let us now suppose that the arrival process is a renewal process with arrival rate $n\mu$ and squared-coefficient-of-variation (SCV), i.e., the variance divided by the square of the mean, of the interarrival times c_a^2 . As in [39], we will consider a diffusion limit for the number of sources as seen at arrival

epochs. Let $N(k)$ be the number of sources seen by the k th arrival. We define the state of the n th system to be

$$\xi_n(k) = \frac{N(k) - n}{\sqrt{n}}$$

and the interpolated process $x_n(t)$ as

$$x_n(t) = \xi_n(k), \quad k \leq t < k + 1.$$

We assume that the rate of departures between successive arrivals remains unchanged. The error due to this approximation is asymptotically negligible [39].

We compute $m(x)$ as

$$\begin{aligned} m(x) &= \lim_{n \rightarrow \infty} E\left(\frac{\xi_n(k+1) - \xi_n(k)}{1/n} \middle| \xi_n(k) = x\right) \\ &= \lim_{n \rightarrow \infty} \sqrt{n}(1 + E(D_n(x\sqrt{n} + n) | a_n)) \end{aligned}$$

where a_n is the interarrival time with mean $1/n\mu$ and $D(\cdot)$ is a Poisson random variable with

$$E(D_n(\sqrt{n}x + n) | a_n) = \begin{cases} \mu(\sqrt{n}x + n)a_n, & \text{if } x < \gamma \\ \mu C_n/\beta_n a_n, & \text{if } x > \gamma. \end{cases}$$

The diffusion coefficient $\sigma^2(x)$ is computed as

$$\begin{aligned} \sigma^2(x) &= \lim_{n \rightarrow \infty} E\left(\frac{(\xi_n(k+1) - \xi_n(k))^2}{1/n} \middle| \xi_n(k) = x\right) \\ &= \lim_{n \rightarrow \infty} E\left(E(1 - D_n(x\sqrt{n} + n | a_n))^2\right). \end{aligned}$$

Carrying out the above calculations yields

$$m(x) = \begin{cases} \theta - \gamma, & \text{if } x > \gamma \\ -x, & \text{if } x < \gamma \end{cases} \quad (24)$$

$$\sigma^2(x) = 1 + c_a^2. \quad (25)$$

Let $p(x; \gamma, \theta, z)$ denote the steady-state density as a function of the parameters γ , θ , and z , where $z = (1 + c_a^2)/2$. Then, by substituting (24) and (25) in (7), we get

$$p(x) = \begin{cases} K_1 e^{-x^2/2z}, & \text{if } x < \gamma \\ K_2 e^{-(\gamma - \theta)x/z}, & \text{if } x > \gamma \end{cases} \quad (26)$$

where

$$K_1 = \frac{1}{\sqrt{2\pi z}(1 - Q(\gamma/\sqrt{z})) + \frac{e^{-\gamma^2/2z}}{(\gamma - \theta)/\sqrt{z}}} \quad (27)$$

and

$$K_2 = \frac{e^{\gamma^2/2z} e^{-\theta\gamma/z}}{\sqrt{2\pi z}(1 - Q(\gamma/\sqrt{z})) + \frac{e^{-\gamma^2/2z}}{(\gamma - \theta)/\sqrt{z}}}. \quad (28)$$

Comparing (26) to (20), it is easy to see that

$$p(x; \gamma, \theta, z) = \frac{1}{\sqrt{z}} p\left(\frac{x}{\sqrt{z}}; \frac{\gamma}{\sqrt{z}}, \frac{\theta}{\sqrt{z}}, 1\right). \quad (29)$$

The weak convergence result for the $GI/M/C$ model is a special case of the $GI/GI/C$ model studied in Section V. However, we consider the two models separately because, in the

$GI/M/C$ case, the stationary distribution of the limiting diffusion process can be explicitly calculated as we have shown, whereas this is not the case for the $GI/GI/C$ model. Unlike the $M/M/C$ model, we do not provide a proof of the fact that the stationary distribution of the n th converges to the stationary distribution of the limiting diffusion process. We do not consider this problem further in this paper, but we simply remark that it may be possible to show this along the lines of a similar result for the $GI/M/C$ FIFO queue in [15]. We summarize the main results of this section in the following theorem. The proof of this theorem is in Appendix III.

Theorem 3: Suppose that $(N(0) - n/\sqrt{n}) \Rightarrow x(0)$; then

$$\frac{N(n\cdot) - n}{\sqrt{n}} \Rightarrow x(\cdot)$$

where $x(t)$ is given the stochastic differential equation (14), $m(x)$ is given by (24), and σ is given by (25). The steady-state distribution of $x(t)$ is given by (26). \diamond

IV. INSENSITIVITY IN THE $M/GI/C$ MODEL WITH MPS

Let us first revisit the $M/M/C$ model. We denote the steady-state probability that there are n sources in the system by π_n . The local balance equation for π_n is given by

$$\lambda\pi_n = \mu r(n+1)\pi_{n+1} \quad (30)$$

where

$$r(n) = \begin{cases} n, & \text{if } n \leq C \\ \beta C, & \text{if } n > C. \end{cases} \quad (31)$$

We want to show that this steady-state distribution is insensitive to the service-time distribution in the $M/GI/C$ model with MPS. We can appeal to the results in [22] to prove this. It is easy to see that our model is a special case of the symmetric queue considered in [22], and thus, the result follows. We also provide a simple, alternate proof here along the lines of the proof in [16]. To this end, define (n, \mathbf{x}) to be the state of the system, where n is the number of sources in the system and \mathbf{x} is an n -dimensional vector of the remaining service times of the n sources. More precisely, x_1 , the first component of \mathbf{x} , is the remaining service time of the first source chosen at random from the n sources, x_2 , the second component of \mathbf{x} is the remaining service time of the second source chosen at random from the remaining $(n-1)$ sources, and so on. Let $p(n, \mathbf{x})$ be the steady-state density of the Markov process, and define $\mathcal{P}\mathbf{x}$ to be the set of all permutations of the vector \mathbf{x} . By considering some time in steady state with the state (n, \mathbf{x}) , and some other time δ units prior to this time, we obtain the following backward equation:

$$\begin{aligned} \sum_{\mathbf{y} \in \mathcal{P}\mathbf{x}} p(n, \mathbf{y}) &= \sum_{\mathbf{y} \in \mathcal{P}(\{\mathbf{x}, 0\})} p(n+1, \mathbf{y})r(n+1)\delta \\ &+ \sum_{\mathbf{y} \in \mathcal{P}(\mathbf{x} \setminus x_i)} p(n-1, \mathbf{y})\lambda\delta f(x_i) \\ &+ \sum_{\mathbf{y} \in \mathcal{P}(\mathbf{x} + \mathbf{1}r(n)\delta)} p(n, \mathbf{y})(1 - \lambda\delta) \end{aligned}$$

where $f(x)$ is the service-time pdf, $\mathbf{x} \setminus x_i$ denotes the $(n-1)$ -dimensional vector obtained by removing the i th component from

\mathbf{x} , $\{\mathbf{x}, 0\}$ denotes the $(n+1)$ -dimensional vector obtained by adding zero to \mathbf{x} , and $\mathbf{1}$ is a vector whose elements are all equal to one. We have ignored $O(\delta^2)$ terms in the above equation. In the limit $\delta \rightarrow 0$, we get

$$\begin{aligned} & - \sum_{i=1}^n \sum_{\mathbf{y} \in \mathcal{P}\mathbf{x}} \frac{\partial p}{\partial y_i} r(n) \\ &= \sum_{\mathbf{y} \in \mathcal{P}(\{\mathbf{x}, 0\})} p(n+1, \mathbf{y})r(n+1) \\ &+ \sum_{i=1}^n \sum_{\mathbf{y} \in \mathcal{P}(\mathbf{x} \setminus x_i)} p(n-1, \mathbf{y})\lambda f(x_i) \\ &- \lambda \sum_{\mathbf{y} \in \mathcal{P}\mathbf{x}} p(n, \mathbf{y}). \end{aligned} \quad (32)$$

We want to show that

$$p(n, \mathbf{x}) = \pi_n \mu^n \prod_{i=1}^n F^c(x_i) \quad (33)$$

where $F^c(x)$ is the complementary cumulative distribution function (ccdf) given by $\int_x^\infty f(s) ds$. The intuition behind the above form for $p(n, \mathbf{x})$ is obtained from a similar result for $M/GI/\infty$ queues in [37]. The form of $p(n, \mathbf{x})$ is equivalent to showing that, conditioned on there being n sources in the system, the remaining-life distribution of each source is the stationary-excess distribution $F_e(x)$ given by

$$F_e(x) = \mu \int_0^x F^c(s) ds. \quad (34)$$

Substituting (33) in (32), using the fact that there are $n!$ permutations of an n -dimensional vector \mathbf{x} , and rearranging terms, we obtain

$$\begin{aligned} & (n-1)! \mu^{n-1} \sum_{i=1}^n (\pi_n r(n)\mu - \lambda\pi_{n-1}) \\ & \cdot \left(\prod_{j=1, j \neq i}^n F^c(x_j) \right) f(x_i) \\ &= n! \mu^n \left(\prod_{j=1}^n F^c(x_j) \right) (\pi_{n+1} r(n+1)\mu - \lambda\pi_n). \end{aligned}$$

Using the balance equation (30), it is easy to see that the above equation is indeed satisfied, thus proving the insensitivity property.

V. $GI/GI/C$ QUEUE WITH MPS

Let us consider a two-phase phase-type service time distribution consisting of two randomly stopped sequence of exponential phases. After spending a mean time of $1/\mu$ in the first phase, the source enters the second phase with probability p , where the mean time spent is again $1/\mu$. Thus, the probability that the customer leaves the system after phase 1 is $(1-p)$. Let

$$n = \frac{\lambda_n(1+p)}{\mu}, \quad C_n = n + \gamma\sqrt{n}, \quad \beta_n = 1 - \theta/\sqrt{n}$$

and define

$$\xi_{1n}(k) = \frac{N_1(k) - n/(1+p)}{\sqrt{n}}$$

and

$$\xi_{2n}(k) = \frac{N_2(k) - pn/(1+p)}{\sqrt{n}}$$

where $N_1(k)$ and $N_2(k)$ are the number of customers in phases 1 and 2, respectively, at the k th arrival epoch. Also, as in Section III, define $x_{1n}(t)$ and $x_{2n}(t)$ as the interpolated versions for $\xi_{1n}(k)$ and $\xi_{2n}(k)$, respectively, and let $x = x_1 + x_2$. Now, as in Section III, we derive the following quantities.

For $x > \gamma$,

$$\begin{aligned} m_1(x_1, x_2) &= \lim_{n \rightarrow \infty} E \left(\frac{\xi_1(k+1) - \xi_1(k)}{1/n} \middle| \xi_1(k) = x_1, \xi_2(k) = x_2 \right) \\ &= \lim_{n \rightarrow \infty} \sqrt{n} \left[1 - \frac{\beta_n \mu C_n N_1(k)}{(N_1(k) + N_2(k)) \lambda_n} \right] \\ &= \lim_{n \rightarrow \infty} \sqrt{n} \\ &\quad \cdot \left[1 - \frac{(n + \gamma \sqrt{n})(1 - \theta/\sqrt{n}) \left(\frac{\sqrt{n} x_1 + n}{(1+p)} \right)}{(\sqrt{n}(x_1 + x_2) + n)n/(1+p)} \right] \\ &= \theta - \gamma - p x_1 + x_2 \end{aligned}$$

$$\begin{aligned} m_2(x_1, x_2) &= \lim_{n \rightarrow \infty} E \left(\frac{\xi_2(k+1) - \xi_2(k)}{1/n} \middle| \xi_1(k) = x_1, \xi_2(k) = x_2 \right) \\ &= \lim_{n \rightarrow \infty} \sqrt{n} \\ &\quad \cdot \left[\frac{-\mu N_2(k) \beta_n C_n}{\lambda_n (N_1(k) + N_2(k))} + \frac{\mu N_1(k) C_n \beta_n p}{\lambda_n (N_1(k) + N_2(k))} \right] \\ &= (1+p)(p x_1 - x_2) \end{aligned}$$

$$\begin{aligned} \sigma_{11}^2(x_1, x_2) &= \lim_{n \rightarrow \infty} \\ &\quad \cdot E \left(\left(\frac{\xi_1(k+1) - \xi_1(k)}{1/n} \right)^2 \middle| \xi_1(k) = x_1, \xi_2(k) = x_2 \right) \\ &= \lim_{n \rightarrow \infty} 1 - \frac{\beta_n C_n N_1(k) \mu E(a_n)}{N_1(k) + N_2(k)} + \frac{\beta_n^2 C_n^2 N_1(k)^2 \mu^2 E(a_n^2)}{(N_1(k) + N_2(k))^2} \\ &= 1 + c_a^2 \end{aligned}$$

$$\begin{aligned} \sigma_{12}^2(x_1, x_2) &= \lim_{n \rightarrow \infty} E \left(\left(\frac{\xi_1(k+1) - \xi_1(k)}{1/n} \right) \left(\frac{\xi_2(k+1) - \xi_2(k)}{1/n} \right) \middle| \right. \\ &\quad \left. \xi_1(k) = x_1, \xi_2(k) = x_2 \right) \\ &= \lim_{n \rightarrow \infty} - \frac{\beta_n C_n N_2(k) \mu E(a_n)}{N_1(k) + N_2(k)} \\ &\quad + \frac{\beta_n^2 C_n^2 N_1(k) N_2(k) \mu^2 E(a_n^2)}{(N_1(k) + N_2(k))^2} - \frac{p \beta_n^2 C_n^2 N_1(k)^2 \mu^2 E(a_n^2)}{(N_1(k) + N_2(k))^2} \\ &= -p \end{aligned}$$

$$\begin{aligned} \sigma_{22}^2(x_1, x_2) &= \lim_{n \rightarrow \infty} E \left(\left(\frac{\xi_2(k+1) - \xi_2(k)}{1/n} \right)^2 \middle| \xi_1(k) = x_1, \xi_2(k) = x_2 \right) \\ &= \lim_{n \rightarrow \infty} \frac{\beta_n^2 C_n^2 N_2(k)^2 \mu^2 E(a_n^2)}{(N_1(k) + N_2(k))^2} + \frac{\beta_n C_n N_2(k) \mu E(a_n)}{N_1(k) + N_2(k)} \\ &\quad + \frac{(p - p^2) \beta_n C_n N_1(k) \mu E(a_n)}{N_1(k) + N_2(k)} \\ &\quad + p^2 \left(\frac{\beta_n \mu C_n N_1(k)}{(N_1(k) + N_2(k)) \lambda_n} + \frac{\beta_n^2 C_n^2 N_1(k)^2 \mu^2 E(a_n^2)}{(N_1(k) + N_2(k))^2} \right) \\ &\quad - \frac{2p \beta_n^2 C_n^2 N_1(k) N_2(k) \mu^2 E(a_n^2)}{(N_1(k) + N_2(k))^2} \\ &= 2p \end{aligned}$$

For $x < \gamma$,

$$\begin{aligned} m_1(x_1, x_2) &= \lim_{n \rightarrow \infty} \sqrt{n} \left[1 - \frac{\mu N_1(k)}{\lambda_n} \right] \\ &= -(1+p)x_1 \\ m_2(x_1, x_2) &= \lim_{n \rightarrow \infty} \sqrt{n} \left[\frac{-\mu N_2(k)}{\lambda_n} + \frac{\mu N_1(k)p}{\lambda_n} \right] \\ &= (1+p)(p x_1 - x_2) \\ \sigma_{11}^2(x_1, x_2) &= \lim_{n \rightarrow \infty} 1 - N_1(k) \mu E(a_n) + N_1^2(k) \mu^2 E(a_n^2) \\ &= 1 + c_a^2 \\ \sigma_{12}^2(x_1, x_2) &= \lim_{n \rightarrow \infty} -N_2(k) \mu E(a_n) \\ &\quad + N_1(k) N_2(k) \mu^2 E(a_n^2) - p N_1^2(k) \mu^2 E(a_n^2) \\ &= -p \\ \sigma_{22}^2(x_1, x_2) &= \lim_{n \rightarrow \infty} N_2(k) \mu E(a_n) + N_2^2(k) \mu^2 E(a_n^2) \\ &\quad + (p - p^2) N_1(k) \mu E(a_n) \\ &\quad + p^2 (\mu N_1(k) + \mu^2 N_1^2(k) E(a_n^2)) \\ &\quad - 2p N_1(k) N_2(k) \mu^2 E(a_n^2) \\ &= 2p. \end{aligned}$$

Let $\mathbf{m}(\mathbf{x})$ denote the drift vector whose i th element is $m_i(\mathbf{x})$, and Σ denotes the matrix whose (i, j) th element is σ_{ij} . For the two-phase phase-type distribution, \mathbf{m} is a two-dimensional vector and Σ is a 2×2 matrix. In general, for a K -phase phase-type distribution, \mathbf{m} will be a K -dimensional vector and Σ will be a $K \times K$ matrix, both of which can be calculated as in the two-dimensional case. The derivation is similar to the two-phase case, but more tedious, and hence is not shown here. Our weak convergence result for the $GI/GI/C$ model is given below, and the proof is presented in Appendix III.

Theorem 4: Consider the n th $GI/GI/C$ system with arrival process $A(n\tilde{t})$ and service times distributed according to a K -phase phase-type distribution, where each phase has an average duration $1/\mu$, and with probability p , a customer in phase k enters phase $k+1$, and with probability $1-p$, it departs from the system. Thus, the overall mean service time is $\sum_{k=1}^K p^{k-1}/\mu$. Here, $A(t)$ is an arrival process with rate 1 and interarrival time SCV c_a^2 . Let $N_{kn}(t)$ denote the number of customers in phase k , $k = 1, 2, \dots, K$, at time t , and let $\mathbf{N}(t)$ denotes the K -dimensional vector whose k th component

is $N_{kn}(t)$. Also, let $N_n(t)$ be the total number of customers in the system at time t , i.e., $N_n(t) = \sum_{k=1}^K N_{kn}(t)$. Let us define

$$x_{kn}(t) = \frac{N_{kn}(t) - np^{k-1} / \sum_{j=1}^K p^{j-1}}{\sqrt{n}}, \quad \forall k = 1, 2, \dots, K$$

let \mathbf{x}_n be the vector whose k th element is x_{kn} , and let $x_n = \sum_{k=1}^K x_{kn}$. We assume that $\mathbf{x}_n(0) \Rightarrow \mathbf{x}(0)$. Let us define $\mathbf{x}(t)$ to be the K -dimensional diffusion process which is the solution of the stochastic differential equation

$$d\mathbf{x} = \mathbf{m}(\mathbf{x}) dt + \Sigma d\mathbf{w}, \quad \mathbf{x}(0) \text{ given}$$

where $\mathbf{w}(t)$ is a K -dimensional Brownian motion. Then,

$$x_n(\cdot) \Rightarrow x(\cdot)$$

where $x(t) = \sum_{k=1}^K x_k(t)$, and $x_k(t)$ is the k th component of $\mathbf{x}(t)$. \diamond

Remark 2: We point out that we have proved the weak convergence of x_n , and not of \mathbf{x}_n . It is an open issue as to whether or not \mathbf{x}_n converges. For our purposes, the convergence of x_n is sufficient. As in the $GI/GI/\infty$ case [39], the limiting process x is not a Markov process. However, unlike the $GI/GI/\infty$ model, it is not even a Gaussian process here.

Remark 3: We note that the above diffusion coefficients and drift vector are different from the $GI/GI/C$ queueing model with FIFO service studied in [15]. In the FIFO case, in addition to the number of customers in each phase that are currently in service, one also has to keep track of the number of waiting customers. Thus, even for a two-phase service-time distribution, the Markov chain has a three-dimensional state vector. In our case, due to the MPS service policy, the dimension of the state space is equal to the number of phases as in the infinite-server model [39]. Further, in the FIFO case, weak convergence to the diffusion limit could not be established in [15], while we can prove this for the MPS policy, as shown in Appendix III. \diamond

Unlike in the $GI/M/C$ case, it does not seem possible to obtain the stationary distribution of this limiting diffusion process. In Section VI, we resort to approximations. However, the fact that this scaling allows us to show weak convergence partially justifies using the approximation.

VI. TRAFFIC ENGINEERING FOR BEST EFFORT SOURCES

We consider the two QoS metrics to evaluate the performance of the congestion-control scheme. We know that the average service time for a source is given by $1/\mu$ if $C = \infty$. If $C < \infty$, let R_C denote the mean response time for a source. Here, response time refers to the total time spent in the system by a source. Then, our performance measure for the congestion-control scheme is the ratio of R_C to $1/\mu$, i.e.,

$$J_C^1 = \mu R_C. \quad (35)$$

The resource allocation or traffic engineering problem is to choose C such that $J_C^1 \leq J^*$, where J^* is some desired performance level. Another performance measure that may be of interest is the probability that the system is in a congested state, i.e., the network, and not the access speed, is the bottleneck:

$$J_C^2 = P(N > C). \quad (36)$$

In the following subsections, we consider these performance measures in the context of the diffusion limits developed earlier.

A. ON-OFF Source Model

As mentioned before, for the ON-OFF source model, the stationary distribution has been proved to be insensitive to the ON and OFF time distributions in [16]. Thus, one can directly compute the performance measures of interest by solving for the stationary distribution of a birth-death Markov chain [16]. However, the diffusion approximation gives remarkably simple rules of thumb for capacity planning, just as in the case of the classical Erlang loss model. For instance, consider the case $\theta = 0$, i.e., the congestion-control scheme is very efficient, and the probability of congestion is the performance measure of interest. Now, suppose that the system is being operated assuming that the number of sources at any instant of time is some fixed amount. If, due to traffic changes, the number of sources increases, then the amount of additional capacity needed to maintain the desired performance can be computed simply by noting that, as long as $\gamma = (S\rho - C/\sqrt{S\rho})$ remains the same, the performance of the system remains the same.

Suppose that response time is the performance measure of interest. Let N be the number of ON sources in steady state. Then, the mean utilization of the system is $E(r_N)$, where r_N was defined in (31). Since the average number of ON sources is $E(N)$, the average rate delivered to each source is $E(r_N)/E(N)$. Thus,

$$J_C^1 = \frac{\mu E(N)}{E(r_N)}.$$

Assuming $\theta = 0$, and using the diffusion approximation, this can be computed approximately as

$$J_C^1 = \frac{\int_{-\infty}^{\infty} (\sqrt{S}x + S\rho)p(x) dx}{\int_{-\infty}^{\gamma\sqrt{\rho}} (\sqrt{S}x + S\rho)p(x) dx + \int_{\gamma\sqrt{\rho}}^{\infty} (S\rho + \gamma\sqrt{S\rho})p(x) dx}$$

For large values of S ,

$$\sqrt{S\rho}(J_C^1 - 1) \approx \int_{\gamma\sqrt{\rho}}^{\infty} \left(\frac{x}{\sqrt{\rho}} - \gamma \right) p(x) dx$$

or, equivalently,

$$J_C^1 \approx 1 + \frac{1}{\sqrt{S\rho}} \int_{\gamma\sqrt{\rho}}^{\infty} \left(\frac{x}{\sqrt{\rho}} - \gamma \right) p(x) dx.$$

Thus, when S increases, if we increase C to maintain a constant value of γ , then J_C^2 remains the same, and J_C^1 decreases, which are both desirable. The above conclusions hold even when $\theta \neq 0$.

B. Renewal Process Model for Source Arrivals

From Little's law,

$$\mu R_c = \frac{\mu E(N)}{\lambda} \approx 1 + \frac{\bar{x}}{\sqrt{n}}$$

where $\bar{x} = E(x)$. For the general $GI/GI/C$ model with MPS, we do not have the limiting distribution of x . However, for the $GI/M/C$ model, we have the exact form of the steady-state distribution of the diffusion limit process given by (26). For the $GI/GI/C$ model, Theorem 4 in Section V shows that the scaling is appropriate for convergence to a well-defined stochastic process in the limit, i.e., $x_n(t) = (N(t) - n/\sqrt{n})$ converges weakly. Thus, we borrow the following heuristic suggested for loss models in [40]. For the $GI/GI/C$ model, we assume that the steady-state distribution satisfies

$$p(x; \gamma, \theta, z) = \frac{1}{\sqrt{z}} p\left(\frac{x}{\sqrt{z}}; \frac{\gamma}{\sqrt{z}}, \frac{\theta}{\sqrt{z}}, 1\right) \quad (37)$$

where z is given by

$$z = 1 + \mu(c_a^2 - 1) \int_0^\infty G^c(t)^2 dt \quad (38)$$

$G(t)$ is the service-time distribution, and $G^c(t) = 1 - G(t)$. In other words, we simply assume the same relationship as in the $GI/M/C$ case, but with z redefined appropriately. The factor z is the heavy-traffic approximation to the *peakedness* parameter commonly used in teletraffic engineering [11]. Peakedness is the ratio of the variance to the mean number of busy servers in a $G/GI/\infty$ model. Unlike the case of loss models, in addition to x , γ and θ also have to be divided by \sqrt{z} . This is partially justified by our results for the $GI/M/C$ case in Section III. Note that, for the Poisson source model, $c_a^2 = 1$, and hence $z = 1$. This is consistent with our result on the insensitivity of steady-state distribution of the $M/GI/C$ queue with MPS to the service-time distribution.

Now, we compute $\bar{x}(\gamma, \theta, z)$ for general interarrival and service-time distributions, and develop a simple formula relating this to the \bar{x} when the arrival process is Poisson and the service times are exponential

$$\begin{aligned} \bar{x}(\gamma, \theta, z) &= \int_{-\infty}^{\infty} xp(x, \gamma, \theta, z) dx \\ &= \frac{1}{\sqrt{z}} \int_{-\infty}^{\infty} xp(x/\sqrt{z}, \gamma/\sqrt{z}, \theta/\sqrt{z}, 1) dx \\ &= \sqrt{z} \int_{-\infty}^{\infty} yp(y, \gamma/\sqrt{z}, \theta/\sqrt{z}, 1) dy \\ &= \sqrt{z} \bar{x}(\gamma/\sqrt{z}, \theta/\sqrt{z}, 1). \end{aligned} \quad (39)$$

Thus,

$$J_C^1 = 1 + \frac{\sqrt{z}}{\sqrt{n}} \bar{x}(\gamma/\sqrt{z}, \theta/\sqrt{z}, 1). \quad (40)$$

Since we have $z = 1$ for the $M/M/C$ case, to compute the performance measure for general renewal arrival processes and

general service-time distributions, we simply compute \bar{x} for an appropriate $M/M/C$ model with MPS, and substitute it in (40) to obtain the desired result.

A similar result also holds for the other performance measure (36), i.e., the probability of congestion:

$$\begin{aligned} J_C^2 &= P(N(\gamma, \theta, z) > C) \\ &\approx P(x(\gamma, \theta, z) > \gamma) \\ &= \int_{\gamma}^{\infty} p(x, \gamma, \theta, z) dx \\ &= \int_{\gamma}^{\infty} \frac{1}{\sqrt{z}} p(x/\sqrt{z}, \gamma/\sqrt{z}, \theta/\sqrt{z}, 1) dx \\ &= \int_{\gamma/\sqrt{z}}^{\infty} p(x, \gamma/\sqrt{z}, \theta/\sqrt{z}, 1) dx \\ &= P(N(\gamma/\sqrt{z}, \theta/\sqrt{z}, 1) > \rho + \gamma\sqrt{\rho}/\sqrt{z}). \end{aligned} \quad (41)$$

Thus, the probability of congestion can be computed for any general interarrival and service-time distributions by solving a corresponding $M/M/C$ model with MPS.

VII. NUMERICAL AND SIMULATION RESULTS

In this section, we present some numerical and simulation results to support the results of the previous sections.

Example 1: We consider ON-OFF sources with exponential ON and OFF times. The pdf and cdf of the scaled number of active sources x are plotted in Figs. 2 and 3 for various values of n , the number of sources. The shapes of the curves in Fig. 2 validate that two Gaussian densities make up the pdf. The curves in Fig. 2 show that, with proper scaling, the cdf's of the number of ON sources are very close to each other, independent of the total number of sources.

Example 2: Figs. 4 and 5 show the pdf and cdf of x for a Poisson source model. The figures lead to a similar conclusion for Poisson sources as Example 1 does for ON-OFF sources.

Example 3: We consider $GI/M/C$ and $GI/GI/C$ models with MPS. The interarrival and service times are taken to be hyperexponential distributions with balanced means, with $c^2 = 1$, corresponding to the exponential distribution. Simulation results are compared with theoretical results in Tables I and II. By theoretical results, we mean the computation of the performance measures using appropriately scaled $M/M/C$ models based on the formulas (40) and (41). The simulations are based on 20 million arrivals divided into 20 batches for each case. The tables show that the simulation and theoretical results agree reasonably well. The results for the $GI/M/C$ case, especially J_C^1 , are better than the corresponding results for the $GI/GI/C$ case. This is to be expected since the approximation for the $GI/GI/C$ is only partially justified, as discussed in Section VI.

VIII. CONCLUSIONS

We have developed diffusion approximations, and have established weak convergence results for congestion-controlled

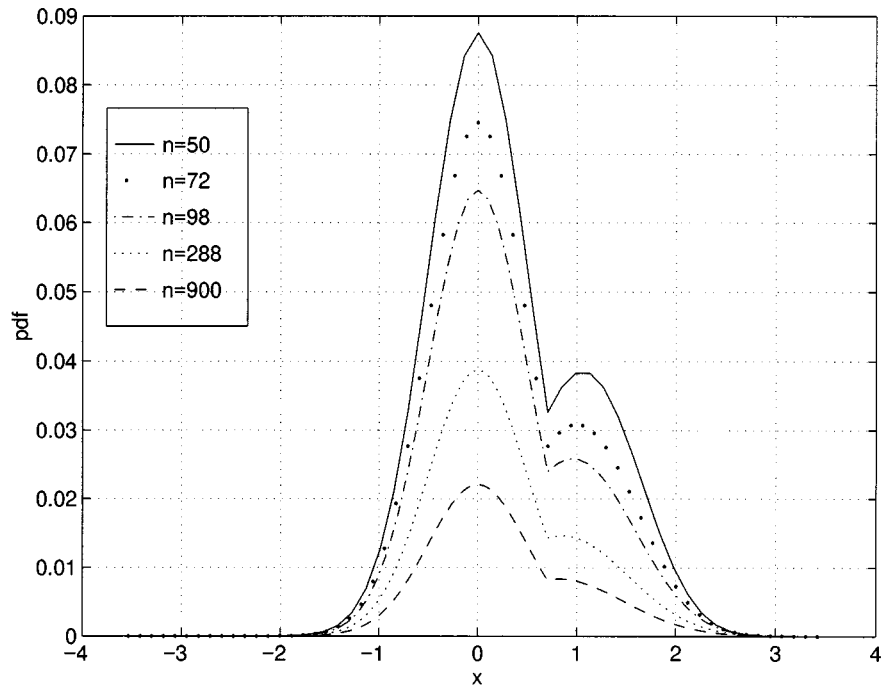


Fig. 2. Example 1: ON-OFF sources pdf, $\lambda = \mu = 1$, $k = 2$, $\gamma = 1$.

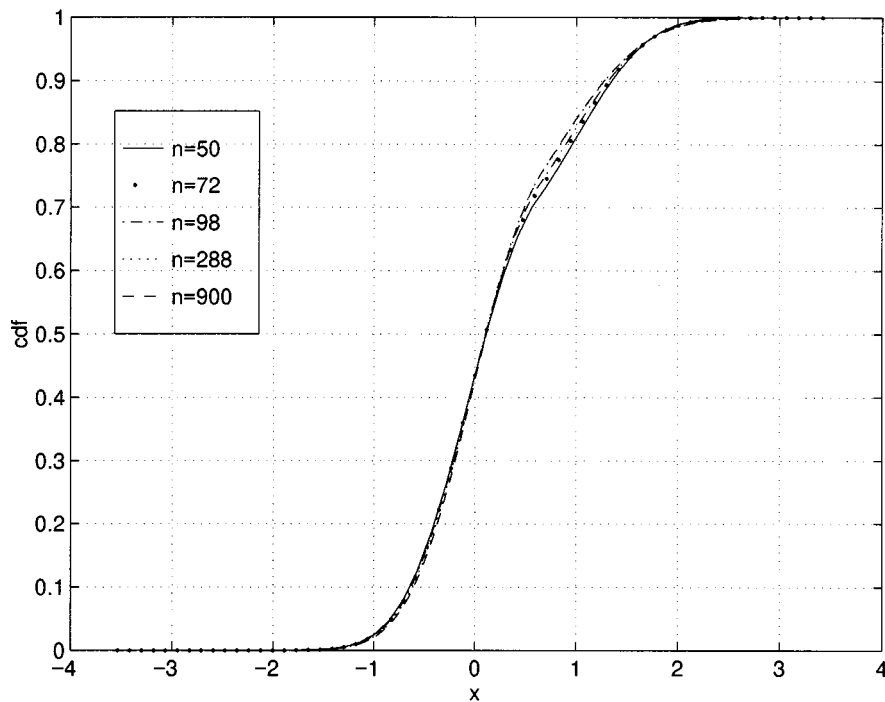


Fig. 3. Example 1: ON-OFF sources cdf, $\lambda = \mu = 1$, $k = 2$, $\gamma = 1$.

queues which serve best effort traffic. Here, congestion control is modeled as a modified processor-sharing policy. The main theoretical result is the weak convergence of an appropriately scaled number-in-the-stem for $GI/GI/C$ models. When the service times are exponential, we are further able to explicitly compute the stationary distribution. These results could be helpful in engineering networks which support best effort traffic,

where the traffic patterns have nonexponential interarrival and service times. Several numerical examples support the use of the approximations.

There are several possible extensions of this work that would merit further study. As in [32] for loss models, it would be interesting to study the applicability of the diffusion approximations in estimating the behavior of the congestion-controlled system

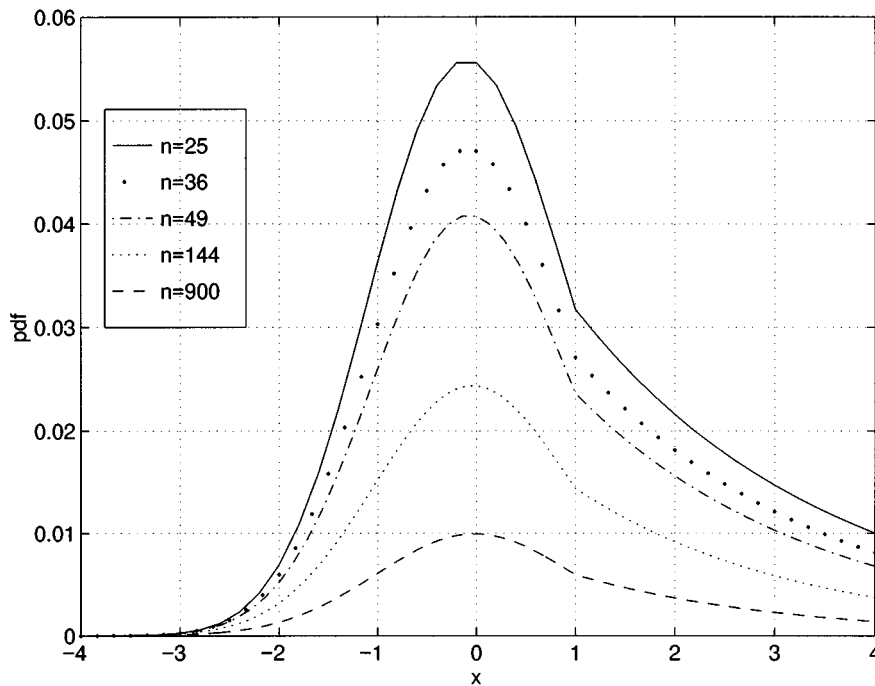


Fig. 4. Example 2: Poisson source pdf, $\mu = 1, k = 0.5, \gamma = 1$.

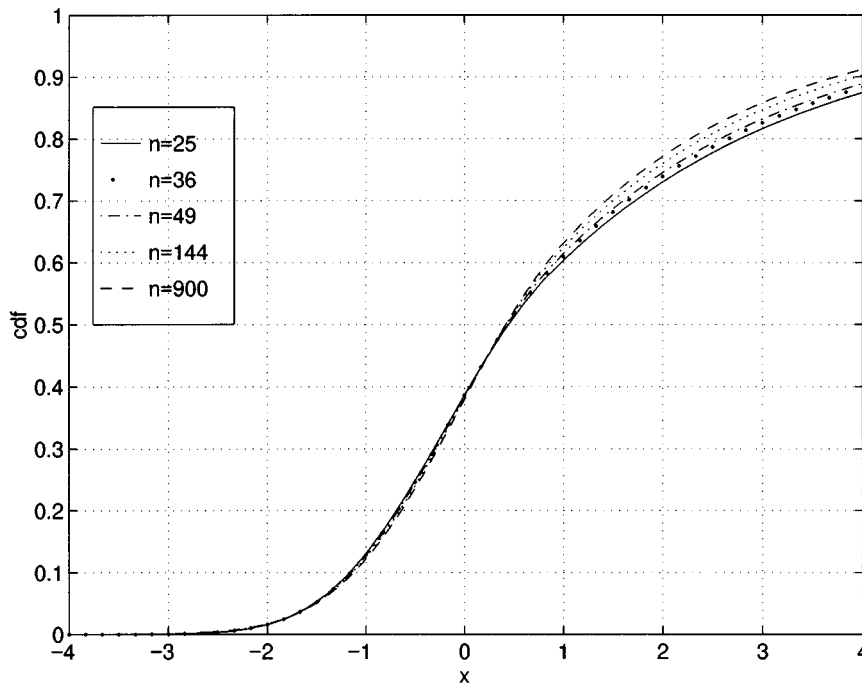


Fig. 5. Example 2: Poisson source cdf, $\mu = 1, k = 0.5, \gamma = 1$.

over finite-time intervals. A second extension would be to study the case where the access link speeds are different for different sources. Finally, numerical results suggest that there are other approximations which could perform better than diffusion approximations [25]. The work in [25] for the $G/M/C$ queue is particularly relevant since it is a special case of the $G/M/C$ model with MPS where $\beta = 0$. It would be interesting to study the extensions of [25], and numerically compare our approxi-

mation to it. In addition, as in [32, Sect. 1.5], one may have to do further refinements based on empirical observations to accurately estimate discrete probabilities using diffusion approximations. If one is interested in estimating discrete probabilities (as opposed to averages or complementary distributions as in this paper), it would be worthwhile to study further refinements to the diffusion approximations or to consider alternate approximations such as [25].

TABLE I
SIMULATED VERSUS THEORETICAL RESULTS FOR $GI/M/C$ CASES

n	c_a^2	c_s^2	z	J_C^1 (sim)	J_C^1 (theory)	J_C^2 (sim)	J_C^2 (theory)
25	10	1	5.5	3.1378	2.8651	0.7470	0.7445
100	10	1	5.5	1.9294	1.8773	0.7329	0.7324
900	10	1	5.5	1.2752	1.2695	0.7307	0.7241
25	5	1	3	2.0542	1.8957	0.6502	0.6559
100	5	1	3	1.4424	1.4080	0.6302	0.6380
900	5	1	3	1.1295	1.1379	0.6153	0.6258

TABLE II
SIMULATED VERSUS THEORETICAL RESULTS FOR $GI/M/C$ CASES

n	c_a^2	c_s^2	z	J_C^1 (sim)	J_C^1 (theory)	J_C^2 (sim)	J_C^2 (theory)
25	5	10	2.1827	1.5243	1.6602	0.5593	0.5992
100	5	10	2.1827	1.2097	1.2722	0.5396	0.5774
900	5	10	2.1827	1.0501	1.0930	0.4938	0.5627
25	10	5	4.0020	2.0045	2.3087	0.6687	0.7011
100	10	5	4.0020	1.4543	1.5785	0.6670	0.6862
900	10	5	4.0020	1.1284	1.1834	0.6380	0.6760

APPENDIX I

LIMIT THEOREM FOR THE PROBABILITY OF CONGESTION IN THE ON-OFF SOURCE MODEL

As seen previously in Section II-A, the MPS system with ON-OFF sources can be modeled as a birth-death process with transition rates q_{ij} given by (1) and (2). We can solve for the steady-state distribution p_i as

$$p_i = \binom{S}{i} \left(\frac{\lambda}{\mu}\right)^i p_0, \quad i < C \quad (42)$$

and

$$p_i = \binom{S}{i} \left(\frac{\lambda}{\mu}\right)^i \frac{(\beta C)^{C-i}}{i!C!} p_0, \quad i \geq C \quad (43)$$

where

$$p_0 = \frac{1}{\left[\sum_{i=0}^{C-1} \binom{S}{i} \left(\frac{\lambda}{\mu}\right)^i \frac{1}{i!} + \sum_{i=C}^n \binom{S}{i} \left(\frac{\lambda}{\mu}\right)^i \frac{(\beta C)^{C-i}}{i!C!} \right]}. \quad (44)$$

Using the diffusion scaling

$$S = n, \quad C = n\rho + \gamma\sqrt{n\rho}$$

$$\beta = 1 - \frac{\theta}{\sqrt{n\rho}}, \quad x(t) = \frac{N(t) - n}{\sqrt{n}}$$

$$\rho = \frac{\lambda}{\lambda + \mu}$$

and

$$\frac{\lambda}{\mu} = \frac{\rho}{(1-\rho)}$$

we want to obtain the expression for $P(N(\infty) > C)$, where $N(t) = \lim_{n \rightarrow \infty} N_n(t)$. Define the following:

$$\alpha_n = P(N(\infty) \geq C)$$

$$= \sum_{i=C}^n \binom{S}{i} \left(\frac{\lambda}{\mu}\right)^i \frac{(\beta C)^{C-i}}{i!C!} p_0 \quad (45)$$

and

$$\alpha = P(x \geq \gamma\sqrt{\rho}). \quad (46)$$

From (11) we can find α , shown in (47) at the bottom of the page.

Our goal in this section is to show that $\lim_{n \rightarrow \infty} \alpha_n = \alpha$. We can rewrite α_n in (45) as

$$\alpha_n = \left[1 + \frac{\vartheta_n^1}{\xi_n^1} \right]^{-1} \quad (48)$$

where

$$\vartheta_n^1 = \sum_{i=0}^{C-1} \binom{S}{i} \rho^i (1-\rho)^{n-i} \quad (49)$$

and

$$\xi_n^1 = (1-\rho)^n \sum_{i=C}^n \binom{S}{i} \left(\frac{\rho}{(1-\rho)}\right)^i \frac{(\beta C)^{C-i}}{i!C!}. \quad (50)$$

Now, we show the desired result by establishing several intermediate facts.

Fact 1:

$$\lim_{n \rightarrow \infty} \vartheta^1 = \left(1 - Q\left(\frac{\gamma}{\sqrt{\rho(1-\rho)}}\right) \right).$$

$$\alpha = \left[1 + \frac{e^{(-1/2(1-\rho))(\gamma\sqrt{\rho} - ((\theta-\gamma)(1-\rho)/\sqrt{\rho}))^2} \sqrt{2\pi\rho(1-\rho)} \left(1 - Q\left(\frac{\gamma\sqrt{\rho}}{\sqrt{\rho(1-\rho)}}\right) \right)}{e^{(-\gamma^2/2\rho(1-\rho))} \sqrt{2\pi(1-\rho)} Q\left(\frac{\gamma - \frac{(\theta-\gamma)(1-\rho)}{\sqrt{\rho}}}{\sqrt{1-\rho}}\right)} \right]^{-1} \quad (47)$$

Proof:

$$\vartheta_n^1 = P(X_n \leq C - 1)$$

where X_n is a binomial random variable with mean = $n\rho$ and variance = $n\rho(1 - \rho)$. Substituting $C = n\rho + \gamma\sqrt{n\rho}$ and using the central limit theorem, we get the desired result. \diamond

Fact 2:

$$\xi_n^1 = \frac{(1 - \rho)^n p_C(J!)}{\left(\frac{1}{\nu}\right)^J e^{-(1/\nu)}} \sum_{i=0}^J \frac{\left(\frac{1}{\nu}\right)^i}{i!} e^{-(1/\nu)}$$

where $\nu = (\rho/\beta C(1 - \rho))$ and $J = n - C$.

Proof:

$$\xi_n^1 = (1 - \rho)^n \sum_{i=C}^n p_i = (1 - \rho)^n \sum_{m=0}^{n-C} p_{C+m}. \quad (51)$$

Note that $p_{C+m} = p_C/K$, where K is given by

$$K = \frac{\binom{n}{C} \left(\frac{\rho}{(1 - \rho)}\right)^C}{\binom{n}{C+m} \left(\frac{\rho}{(1 - \rho)}\right)^{C+m} \frac{(\beta C)^{C-(C+m)} (C+m)!}{C!}} = \frac{\left(\frac{\beta C(1 - \rho)}{\rho}\right)^m}{\binom{n-C}{m} m!}.$$

Substituting in (51) proves the fact. \diamond

Next, we rewrite ξ_n^1 as

$$\xi_n^1 = \xi_n^2 \xi_n^3$$

where

$$\xi_n^2 = \frac{(1 - \rho)^n p_C(J!)}{\left(\frac{1}{\nu}\right)^J e^{-(1/\nu)}} \quad \text{and} \quad \xi_n^3 = \sum_{i=0}^J \frac{\left(\frac{1}{\nu}\right)^i}{i!} e^{-(1/\nu)}.$$

Fact 3:

$$\lim_{n \rightarrow \infty} \xi_n^3 = Q\left(\frac{\gamma\sqrt{\rho} - (\theta - \gamma)\frac{(1 - \rho)}{\sqrt{\rho}}}{\sqrt{(1 - \rho)}}\right).$$

Proof:

$$\xi_n^3 = \sum_{i=0}^J \frac{\left(\frac{1}{\nu}\right)^i}{i!} e^{-(1/\nu)} = P(X_n \leq J)$$

where X_n is a Poisson random variable with mean = variance = $1/\nu$. Recall $n - C = J$ and $\nu = (\rho/\beta C(1 - \rho))$. Using the central limit theorem, we get the desired result by using the fact that

$$\lim_{n \rightarrow \infty} \frac{n - C - \frac{\beta C(1 - \rho)}{\rho}}{\sqrt{\frac{\beta C(1 - \rho)}{\rho}}} = \frac{-\gamma\sqrt{\rho} + (\theta - \gamma)\frac{(1 - \rho)}{\sqrt{\rho}}}{\sqrt{(1 - \rho)}}.$$

\diamond

Let $\xi_n^2 = \xi_n^4 \xi_n^5$, where

$$\xi_n^4 = (1 - \rho)^n p_C \quad \xi_n^5 = \frac{J!}{\left(\frac{1}{\nu}\right)^J} e^{-(1/\nu)}.$$

Fact 4:

$$\lim_{n \rightarrow \infty} \sqrt{n\rho(1 - \rho)} \xi_n^4 = \frac{e^{-(\gamma\sqrt{\rho})^2/2\rho(1 - \rho)}}{\sqrt{2\pi}}.$$

Proof:

$$\begin{aligned} \xi_n^4 &= (1 - \rho)^n p_C = \binom{n}{C} \rho^C (1 - \rho)^{n-C} \\ &= P\left(C < \sum_{i=1}^n Y_{i=1} \leq C + 1\right) \end{aligned}$$

where Y_i are i.i.d. Bernoulli random variables with mean = $n\rho$ and variance = $n\rho(1 - \rho)$. Thus, as shown in the equation at the bottom of the page, where the last line above follows from the central limit theorem. \diamond

Fact 5:

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{1}{\sqrt{2\pi(n - C)}} \xi_n^5 \\ = e^{(1/2(1 - \rho))(\gamma\sqrt{\rho} - (\theta - \gamma)((1 - \rho)/\sqrt{\rho}))^2}. \end{aligned}$$

Proof: Applying Stirling's formula, $J! \approx \sqrt{2\pi J} J^J e^{-J}$, and letting $1/\nu = \tau$,

$$\begin{aligned} \xi_n^5 &= \sqrt{2\pi J} \left(\frac{J}{\tau}\right)^J e^{-(J - \tau)} \\ &= \sqrt{2\pi J} e^{J \log(J/\tau) - (J - \tau)}. \end{aligned}$$

$$\frac{1}{\sqrt{n\rho(1 - \rho)}} \xi_n^4 = \frac{P\left(\frac{\gamma\sqrt{\rho}}{\sqrt{\rho(1 - \rho)}} < \frac{\sum_i Y_i - n\rho}{\sqrt{n\rho(1 - \rho)}} \leq \frac{\gamma\sqrt{\rho}}{\sqrt{\rho(1 - \rho)}} + \frac{1}{\sqrt{n\rho(1 - \rho)}}\right)}{\frac{1}{\sqrt{n\rho(1 - \rho)}}} \rightarrow \frac{e^{-(\gamma\sqrt{\rho})^2/2\rho(1 - \rho)}}{\sqrt{2\pi}}$$

Noting the fact that $\log(x) = \log(1 - (1 - x)) = -(1 - x) - (1/2)(1 - x)^2 - o(1 - x)^2$,

$$\begin{aligned}\xi_n^5 &= \sqrt{2\pi J} e^{(J-\tau)+(J/2)(J-\tau/J)^2-(J-\tau)} \\ &= \sqrt{2\pi J} e^{(\tau/2J)(J-\tau\sqrt{\tau})}.\end{aligned}$$

Recalling $J = n - C$ and $\tau = (\beta C(1 - \rho)/\rho)$, and using the fact that

$$\begin{aligned}\lim_{n \rightarrow \infty} \frac{\beta C(1 - \rho)}{2(n - C)} \left(\frac{n - C - \frac{\beta C(1 - \rho)}{\rho}}{\sqrt{\frac{\beta C(1 - \rho)}{\rho}}} \right)^2 \\ = \frac{1}{2} \left(\frac{-\gamma\sqrt{\rho} - (\theta - \gamma) \frac{(1 - \rho)}{\sqrt{\rho}}}{(1 - \rho)} \right)^2\end{aligned}$$

proves this fact. \diamond

The following result is an immediate consequence of Facts 3–5.

Fact 6:

$$\begin{aligned}\lim_{n \rightarrow \infty} \xi_n^1 \\ = \frac{\sqrt{2\pi(1 - \rho)} e^{(1/2(1 - \rho)(\gamma\sqrt{\rho} - (\theta - \gamma)((1 - \rho)/\sqrt{\rho}))^2)}}{\sqrt{2\pi\rho(1 - \rho)} e^{(\gamma\sqrt{\rho})^2/2\rho(1 - \rho)}} \\ \times Q \left(\frac{\gamma\sqrt{\rho} - (\theta - \gamma) \frac{(1 - \rho)}{\sqrt{\rho}}}{\sqrt{(1 - \rho)}} \right).\end{aligned}$$

Using Facts 1 and 6 in (48), we get the expression for α given in (47). \diamond

APPENDIX II

LIMIT THEOREM FOR THE PROBABILITY OF CONGESTION IN THE POISSON SOURCE MODEL

As seen previously in Section II-B, the MPS system with Poisson sources can be modeled as a birth–death process with transition rates q_{ij} given by (15) and (16). We can solve for the steady-state distribution p_i as

$$p_i = \left(\frac{\lambda}{\mu} \right)^i \frac{p_0}{i!}, \quad i < C \quad (52)$$

and

$$p_i = \left(\frac{\lambda}{\beta C \mu} \right)^i \frac{(\beta C)^C}{C!} p_0, \quad i \geq C \quad (53)$$

where

$$p_0 = \frac{1}{\left[\sum_{i=0}^{C-1} \left(\frac{\lambda}{\mu} \right)^i \frac{1}{i!} + \left(\frac{\lambda}{\mu} \right)^C \frac{1}{C! \left(1 - \frac{\lambda}{\beta C \mu} \right)} \right]}. \quad (54)$$

Using the diffusion scaling

$$\begin{aligned}\lambda &= n\mu, \quad C = n + \gamma\sqrt{n} \\ \beta &= 1 - \frac{\theta}{\sqrt{n}}\end{aligned}$$

and

$$x(t) = \frac{N(t) - n}{\sqrt{n}}$$

define the following:

$$\begin{aligned}\alpha_n &= P(N(\infty) \geq C) = \sum_{i=C}^{\infty} \left(\frac{\lambda}{\beta C \mu} \right)^i \frac{(\beta C)^C}{C!} p_0 \\ &= \frac{n^C p_0}{C! \left(1 - \frac{n}{\beta C} \right)}\end{aligned} \quad (55)$$

and

$$\alpha = P(x \geq \gamma). \quad (56)$$

From (20) in Section II-B, we find that

$$\alpha = \frac{e^{(-\gamma^2/2)}}{\sqrt{2\pi}(\gamma - \theta)(1 - Q(\gamma)) + e^{(-\gamma^2/2)}}. \quad (57)$$

Our goal in this section is show that $\lim_{n \rightarrow \infty} \alpha_n = \alpha$. We can rewrite α_n in (55) as

$$\alpha_n = \left[1 + \frac{\vartheta_n}{\xi_n} \right]^{-1} \quad (58)$$

where

$$\vartheta_n = \sum_{i=0}^{C-1} \frac{n^i}{i!} e^{-n}$$

and

$$\xi_n = \frac{n^C e^{-n}}{C! \left(1 - \frac{n}{\beta C} \right)}.$$

Now, we show the desired result by establishing the following intermediate facts.

Fact 7:

$$\lim_{n \rightarrow \infty} \vartheta_n = 1 - Q(\gamma).$$

Proof:

$$\vartheta_n = P(X_n \leq C - 1)$$

where X_n is a Poisson random variable with mean and variance both equal to n . Substituting $C = n\rho + \gamma\sqrt{n\rho}$ and applying the central limit theorem, we prove the fact. \diamond

Fact 8:

$$\lim_{n \rightarrow \infty} \xi_n = \frac{e^{-\gamma^2/2}}{\sqrt{2\pi}(\gamma - \theta)}.$$

Proof: Applying Stirlings formula to ξ_n yields

$$\xi_n = \frac{e^{C(1-(n/C)+\log(n/C))}}{\sqrt{2\pi C} \frac{\beta C - n}{\beta C}}.$$

Using

$$\log(x) = -(1-x) - \frac{(1-x)^2}{2} + o(1-x)^2$$

we have

$$\xi_n = \frac{1}{\sqrt{2\pi}} \frac{e^{(-nC^2(1-(n/C))^2/2Cn)}}{\sqrt{C} \frac{\beta C - n}{\beta C}}.$$

Finally, substituting $C = n + \gamma\sqrt{n}$, and using the facts that

$$\lim_{n \rightarrow \infty} \frac{n}{C} = 1 \quad \text{and} \quad \lim_{n \rightarrow \infty} \sqrt{C} \frac{\beta C - n}{\beta C} = (\gamma - \theta)$$

yields the fact. \diamond

Using Facts 7 and 8 in (58), we get the expression for α in (57).

APPENDIX III PROOFS OF THEOREMS 1, 3, AND 4

We first give the proof for the $GI/M/C$ case (Theorem 3), and later point out the modifications necessary to complete the proof for the $GI/GI/C$ model (Theorem 4) and the ON-OFF model (Theorem 1). The basic idea of the proof is as follows. Consider the scaled, centered, and interpolated process sampled at arrival epochs, denoted by $x_n(t)$. We first construct two processes $x_{n,\epsilon}^u(t)$ and $x_{n,\epsilon}^l(t)$, which upper and lower bound $x_n(t)$, respectively. We use the results in [12, ch. 7, Corollary 4.2, p. 355] to show that $x_{n,\epsilon}^l(t) \Rightarrow x_\epsilon^l(t)$, and $x_{n,\epsilon}^u(t) \Rightarrow x_\epsilon^u(t)$ as $n \rightarrow \infty$. Then, using the results in [36] on the convergence of diffusions, we show that both $x_\epsilon^l(t)$ and $x_\epsilon^u(t)$ weakly converge to $x(t)$. Thus, we would have proved that $x_n(t) \Rightarrow x(t)$. To do this, we need the limiting martingale problem to be well posed. Standard results show that a process $x(t)$, defined by a stochastic differential equation with a discontinuous drift as in our problem, has a weak solution [20, Proposition 3.6, p. 303], and thus, the martingale problem is well posed [20, Proposition 4.11, p. 318]. For a discussion of weak convergence in another context where there is a discontinuity in the state dependence, see [29, Remark, p. 277].

Let us define $N_{n,\epsilon}^l(t)$ as the number of sources at time t in the system whose departure process is given by $A(\int_0^t r_{n,\epsilon}^l(N(s)) ds)$, where $A(t)$ is a unit Poisson process and $r_{n,\epsilon}^l(N)$ is given by

$$r_{n,\epsilon}^l(N) = \begin{cases} \mu N, & N \leq C \\ \beta \mu C, & N > C + \epsilon\sqrt{n} \\ \mu C \left(\frac{(N-C)(1-\beta)}{C-n-(\gamma+\epsilon)\sqrt{n}} + 1 \right), & C < N \leq C + \epsilon\sqrt{n}. \end{cases}$$

Similarly, define $N_{n,\epsilon}^u(t)$ by defining the departure rate $r_{n,\epsilon}^u(N)$ as

$$r_{n,\epsilon}^u(N) = \begin{cases} \mu N, & N \leq C - \epsilon\sqrt{n} \\ \beta \mu C, & N > C \\ (N-C)\mu \left(\frac{\beta C - n - (\gamma-\epsilon)\sqrt{n}}{C-n-(\gamma-\epsilon)\sqrt{n}} \right) + \beta \mu C, & C - \epsilon\sqrt{n} < N \leq C. \end{cases}$$

In the above construction, we simply increase the departure rate when the number of busy servers is in the interval $(C, C + \epsilon\sqrt{n}]$ to obtain $N_{n,\epsilon}^l(t)$, and we decrease the departure rate in $(C - \epsilon\sqrt{n}, C]$ to obtain $N_{n,\epsilon}^u(t)$. Thus, it is easy to see that, for all t and ϵ ,

$$N_{n,\epsilon}^l(t) \leq N_n(t) \leq N_{n,\epsilon}^u(t), \quad \text{w.p. 1.}$$

The above bounds can be rigorously obtained by coupling the three processes appropriately as follows. At time $t = 0$, let the above inequality hold w.p. 1. Let all three processes have the same arrival process sample paths. Then, make the downward jumps of $N_n(t)$ a subset of the downward jumps of $N_{n,\epsilon}^l(t)$, and further, make the downward jumps of $N_{n,\epsilon}^u(t)$ a subset of the downward jumps of $N_n(t)$, which gives us the desired ordering.

As before, define $\xi_{n,\epsilon}^l(k)$ and $\xi_{n,\epsilon}^u(k)$ as

$$\xi_{n,\epsilon}^l(k) = \frac{N_{n,\epsilon}^l(k) - n}{\sqrt{n}} \quad \text{and} \quad \xi_{n,\epsilon}^u(k) = \frac{N_{n,\epsilon}^u(k) - n}{\sqrt{n}}$$

respectively, where $N_{n,\epsilon}^l(k)$ and $N_{n,\epsilon}^u(k)$ are the number of sources seen by the k th arrival in the respective systems. Let $x_{n,\epsilon}^l(t)$ and $x_{n,\epsilon}^u(t)$ be the interpolated processes:

$$x_{n,\epsilon}^l(t) = \xi_{n,\epsilon}^l(\lfloor nt \rfloor) \quad \text{and} \quad x_{n,\epsilon}^u(t) = \xi_{n,\epsilon}^u(\lfloor nt \rfloor).$$

By our construction of $N_{n,\epsilon}^l(t)$ and $N_{n,\epsilon}^u(t)$,

$$x_{n,\epsilon}^u(t) \leq x_n(t) \leq x_{n,\epsilon}^l(t), \quad \text{w.p. 1.}$$

As in Section V, we consider the imbedded Markov chain at arrival epochs, and derive the drift and diffusion coefficients. From [39, Sect. 3.2], the error due to the assumption that the number of sources in the system does not change between arrival instants is asymptotically negligible, and further, the difference between the process at arrival instants and any instant converges to zero. Thus, using [12, ch. 7, Corollary 4.2, p. 355], it follows that $x_{n,\epsilon}^l \Rightarrow x_\epsilon^l$ and $x_{n,\epsilon}^u \Rightarrow x_\epsilon^u$, where the process $x_\epsilon^l(t)$ has the drift vector

$$m_\epsilon^l(x) = \begin{cases} -x, & x \leq \gamma \\ \frac{(x-\gamma)\theta}{\epsilon} - \gamma, & \gamma < x \leq \gamma + \epsilon \\ \theta - \gamma, & x > \gamma + \epsilon \end{cases}$$

and the process $x_\epsilon^u(t)$ has the drift vector

$$m_\epsilon^u(x) = \begin{cases} -x, & x \leq \gamma - \epsilon \\ -(x-\gamma)\frac{\mu}{\epsilon}(\epsilon-\theta) - (\gamma-\theta)\mu, & \gamma - \epsilon < x \leq \gamma \\ \theta - \gamma, & x > \gamma. \end{cases}$$

The above construction is necessary since one cannot apply [12, ch. 7, Corollary 4.2, p. 355] if the limiting diffusion process has

a discontinuous drift. By our construction, $m_\epsilon^l(x)$ and $m_\epsilon^u(x)$ are continuous functions of x . The diffusion coefficient remains the same as that for $x(t)$, i.e., $\sigma^2(x) = 2\mu$.

Now, from [36, Theorem 11.3.4, p. 278], upon verification of some conditions [36, eqs. (3.4)–(3.8), p. 274], it follows that x_ϵ^l and $x_{n,\epsilon}^u$ weakly converge to the same process $x(t)$ in the limit as $\epsilon \rightarrow 0$, thus completing the proof for the $GI/M/C$ case. For the lower bound, the condition [36, eq. (3.4), p. 274] requires that

$$\begin{aligned} \lim_{\epsilon \rightarrow 0} \int_0^\infty \int_0^\infty m_\epsilon^l(x) \phi(s, x) ds dx \\ = \int_0^\infty \int_0^\infty m(x) \phi(s, x) ds dx \end{aligned} \quad (59)$$

$\forall \phi(s, x) \in C_0^\infty([0, \infty) \times \mathcal{R})$, where $C_0^\infty([0, \infty) \times \mathcal{R})$ is the set of all functions with compact support, and which possess continuous derivatives of all orders. Even though there is a division by ϵ (which goes to zero) in the definition of $m_\epsilon^l(x)$, it is easy to check that

$$|m_\epsilon^l(x)| \leq |x|$$

$\forall \epsilon > 0$, and thus, by the Lebesgue dominated convergence theorem [8, p. 213], condition (59) is verified. A similar verification can be done for the upper bound also. The rest of the conditions of [36, Theorem 11.3.4, p. 278] are easy to verify.

Finally, to prove the weak convergence for the $GI/GI/C$ model, we simply outline the steps, and the rest is similar to the $GI/M/C$ case. If there are K phases in the service-time distribution, as in Section V, we would have a K -dimensional Markov chain at arrival instants with state vector

$$\mathbf{N}_n = (N_{1n}, N_{2n}, \dots, N_{Kn}).$$

As before, we define new departure rates, one larger and one smaller, to upper and lower bound $\sum_{i=1}^K N_{in}$. Note that we do not have element-by-element bounds on \mathbf{N}_n ; rather, we upper and lower bound the total number of customers in the system as follows. We allow the original system, and the upper bounding and lower bounding systems to have the same arrival process sample paths. In addition, at each arrival instant, we populate each system with a customer whose service time is drawn from the K -phase phase-type distribution. Let the state of the upper and lower bounding Markov chains be $\mathbf{N}_{n,\epsilon}^u$ and $\mathbf{N}_{n,\epsilon}^l$, respectively. Then, it is easy to show that, $\forall t > 0$,

$$\sum_{i=1}^K N_{1,n,\epsilon}^l(t) \leq \sum_{i=1}^K N_{in}(t) \leq \sum_{i=1}^K N_{1,n,\epsilon}^u(t)$$

w.p.1 if the initial conditions are so ordered, and the remaining service times of at least l customers, where $l = \sum_{i=1}^K N_{1,n,\epsilon}^l(0)$, are the same w.p.1 in all three systems. Next, as in the $GI/M/C$ case, we can show the convergence of centered and scaled versions of $\mathbf{N}_{n,\epsilon}^u$ and $\mathbf{N}_{n,\epsilon}^l$, denoted by \mathbf{X}_ϵ^u and \mathbf{X}_ϵ^l , respectively, to the process \mathbf{x} defined in Theorem 4, where the limit is first taken as $n \rightarrow \infty$, and then $\epsilon \rightarrow 0$, as before. Since summing the elements of a vector-valued random

process is a continuous mapping, we have the desired weak convergence result from the continuous mapping theorem [7].

Finally, for the ON-OFF model, define the upper and lower bounding Markov chains by defining new departure rates as before. Here, instead of coupling arguments, we directly observe that there is a stochastic ordering relationship among the three stochastic processes by checking the simple conditions in [35], Proposition 4.2.10, 2^0 for birth-death Markov chains. The rest of the proof then follows as before.

ACKNOWLEDGMENT

The authors thank W. Whitt, AT&T Laboratories, for many useful discussions.

REFERENCES

- [1] E. Altman, T. Başar, and R. Srikant, "Multi-user rate-based flow control with action delays: A team-theoretic approach," in *Proc. 36th IEEE Conf. Decision Contr.*, San Diego, CA, Dec. 1997.
- [2] —, "Robust rate control for ABR sources," in *Proc. IEEE INFOCOM*, San Francisco, CA, Mar. 1998.
- [3] P. Barford and M. E. Crovella, "Generating representative web workloads for network and server performance evaluation," in *Proc. Performance '98/ACM SIGMETRICS'98*, Madison, WI, 1998, pp. 151–160.
- [4] L. Benmohamed and S. M. Meerkov, "Feedback control of congestion in packet switching networks: The case of a single congested node," *IEEE/ACM Trans. Networking*, vol. 1, no. 6, pp. 693–707, 1993.
- [5] L. Benmohamed and Y. T. Wang, "A control-theoretic ABR explicit rate algorithm for ATM switches with per-VC queueing," in *Proc. IEEE INFOCOM*, San Francisco, CA, Mar. 1998.
- [6] D. Bertsekas and R. Gallager, *Data Networks*. Englewood Cliffs, NJ: Prentice-Hall, 1987.
- [7] P. Billingsley, *Convergence of Probability Measures*. New York: Wiley, 1968.
- [8] —, *Probability and Measure*, 2nd ed. New York: Wiley, 1986.
- [9] A. A. Borovkov, "On limit laws for service processes in multi-channel systems," *Siberian Math. J.*, pp. 746–763. (English transl.).
- [10] —, *Stochastic Processes in Queueing Theory*. New York: Springer-Verlag, 1976.
- [11] A. E. Eckberg, "Generalized peakedness of teletraffic processes," in *Proc. 10th Int. Teletraffic Congr.*, Montreal, Canada, June 1983, p. 4.4b3.
- [12] S. N. Ethier and T. G. Kurtz, *Markov Processes: Characterization and Convergence*. New York: Wiley, 1994.
- [13] P. W. Glynn, "On the Markov property of the $GI/G/\infty$ Gaussian limit," *Adv. Appl. Prob.*, vol. 14, pp. 191–194, 1982.
- [14] P. W. Glynn and W. Whitt, "A new view of the heavy-traffic limit theorem for infinite-server queues," *Adv. Appl. Prob.*, vol. 23, pp. 188–209, 1991.
- [15] S. Halfin and W. Whitt, "Heavy-traffic limits for queues with many exponential servers," *Oper. Res.*, vol. 29, no. 3, pp. 567–588, 1981.
- [16] D. Heyman, T. V. Lakshman, and A. Niedhart, "A new method for analyzing feedback-based protocols with applications to engineering Web traffic over the Internet," in *Sigmetrics 97*, 1997, pp. 24–38.
- [17] D. Iglehart, "Limit diffusion approximations for the many server queue and the repairman problem," *J. Appl. Prob.*, pp. 429–441, 1965.
- [18] V. Jacobson, "Congestion avoidance and control," *ACM Comput. Commun. Rev.*, vol. 18, pp. 314–329, Aug. 1988.
- [19] S. Kalyanaraman, R. Jain, S. Fahmy, R. Goyal, and B. Vandalore. (1997) The ERICA switch algorithm for ABR traffic management in ATM networks. [Online] <http://www.cis.ohio-state.edu/jain/papers>.
- [20] I. Karatzas and S. Shreve, *Brownian Motion and Stochastic Calculus*. New York: Springer, 1996.
- [21] S. Karlin and H. Taylor, *A Second Course in Stochastic Processes*. New York: Academic, 1981.
- [22] F. P. Kelly, *Reversibility and Stochastic Networks*. New York: Wiley, 1976.
- [23] S. Keshav, *An Engineering Approach to Computer Networks*. Reading, MA: Addison-Wesley, 1997.
- [24] L. Kleinrock, *Queueing Systems, Volume 2: Computer Applications*. New York: Wiley-Interscience, 1975.
- [25] C. Knessl, "The WKB approximation to the G/M/m queue," *SIAM J. Appl. Math.*, vol. 51, pp. 1119–1133, 1991.

- [26] H. J. Kushner, "Control of trunk line systems in heavy traffic," *SIAM J. Contr. Optimiz.*, vol. 33, pp. 765–803, 1995.
- [27] —, "Heavy traffic analysis of controlled multiplexing systems," *Queueing Syst.*, vol. 28, pp. 79–107, 1998.
- [28] H. J. Kushner and P. G. Dupuis, *Numerical Methods for Stochastic Control Problems in Continuous Time*. New York: Springer-Verlag, 1992.
- [29] H. J. Kushner and L. F. Martins, "Limit theorems for pathwise average cost per unit time problems for controlled queues in heavy traffic," *Stoch. and Stoch. Rep.*, vol. 42, pp. 25–51, 1993.
- [30] S. Mascolo, D. Cavendish, and M. Gerla, "ATM rate based congestion control using a Smith predictor: An EPRCA implementation," in *Proc. IEEE INFOCOM*, San Francisco, CA, Mar. 1996.
- [31] K. W. Ross, *Multiservice Loss Models for Broadband Telecommunication Networks*. London: Springer, 1995.
- [32] R. Srikant and W. Whitt, "Simulation run lengths to estimate blocking probabilities in multiserver loss models," *ACM Trans. Modelling and Comput. Simulation*, pp. 7–52, Jan. 1996.
- [33] —, "Variance reduction in simulation of loss models," *Oper. Res.*, June 1999.
- [34] C. Stone, "Limit theorems for random walks, birth and death processes and diffusion processes," *Illinois J. Math.*, pp. 638–660, 1961.
- [35] D. Stoyan, *Comparison Methods for Queues and Other Stochastic Models*. New York: Wiley, 1983.
- [36] D. W. Stroock and S. R. S. Varadhan, *Multidimensional Diffusion Processes*. New York: Springer-Verlag, 1979.
- [37] L. Takacs, *Introduction to the Theory of Queues*. Oxford, England: Oxford Univ. Press, 1962.
- [38] P. Varaiya and J. Walrand, *High-Performance Communication Networks*. San Francisco, CA: Morgan Kaufman, 1996.
- [39] W. Whitt, "On the heavy-traffic limit theorems for $GI/G/\infty$ queues," *Adv. Appl. Prob.*, vol. 14, pp. 171–190, 1982.
- [40] —, "Heavy traffic approximations for service systems with blocking," *AT&T Tech. J.*, pp. 689–707, May–June 1984.



Atanu Das received the B.S. and M.S. degrees in electrical engineering from the University of Illinois at Urbana–Champaign in 1996 and 1998, respectively.

From 1996 to 1998, he was a Teaching Assistant at the University of Illinois. He is currently employed as a Systems Engineer with Tellabs Inc., Bolingbrook, IL. His research interests include the modeling and simulation of communication networks.



R. Srikant received the B.Tech. degree from the Indian Institute of Technology, Madras, in 1985, and the M.S. and Ph.D. degrees from the University of Illinois in 1988 and 1991, respectively, all in electrical engineering.

He was a Member of Technical Staff at AT&T Bell Laboratories from 1991 to 1995. Since August 1995, he has been with the University of Illinois, where he is currently an Assistant Professor in the Department of General Engineering and Coordinated Science Laboratory, and an Affiliate in the Department of Electrical and Computer Engineering. His research interests include communication networks, queueing theory, and stochastic control.

Dr. Srikant received an NSF CAREER award in 1997.