

$M_t/G/\infty$ Queues with Sinusoidal Arrival Rates

Stephen G. Eick • William A. Massey • Ward Whitt
AT&T Bell Laboratories, Murray Hill, New Jersey 07974-0636

In this paper we describe the mean number of busy servers as a function of time in an $M_t/G/\infty$ queue (having a nonhomogeneous Poisson arrival process) with a sinusoidal arrival rate function. For an $M_t/G/\infty$ model with appropriate initial conditions, it is known that the number of busy servers at time t has a Poisson distribution for each t , so that the full distribution is characterized by its mean. Our formulas show how the peak congestion lags behind the peak arrival rate and how much less is the range of congestion than the range of offered load. The simple formulas can also be regarded as consequences of linear system theory, because the mean function can be regarded as the image of a linear operator applied to the arrival rate function. We also investigate the quality of various approximations for the mean number of busy servers such as the pointwise stationary approximation and several polynomial approximations. Finally, we apply the results for sinusoidal arrival rate functions to treat general periodic arrival rate functions using Fourier series. These results are intended to provide a better understanding of the behavior of the $M_t/G/\infty$ model and related $M_t/G/s/r$ models where some customers are lost or delayed.

(Queues; Nonstationary Queues; Infinite-server Queues; Queues with Time-dependent Arrival Rates; Approximations; Pointwise Stationary Approximation)

1. Introduction

This paper is part of an effort to analyze queues with time-dependent arrival rates. We want to develop computational methods and approximation techniques, but *the primary purpose of this paper is to develop a better understanding of the time-dependent behavior*. We want to understand the congestion response to different arrival rate functions. For example, we want to understand how peak congestion lags behind the peak arrival rate.

We believe that a good place to begin a quest for better understanding is the $M_t/G/\infty$ queue, which has a nonhomogeneous Poisson arrival process with deterministic time-dependent arrival-rate function $\lambda \equiv \{\lambda(t) : -\infty < t < \infty\}$, i.i.d. service times that are independent of the arrival process and infinitely many servers. *The $M_t/G/\infty$ model is a good starting point because it is remarkably easy to analyze*. For the $M_t/G/\infty$ model we can carry out the mathematical analysis required to obtain explicit formulas. These results apply directly to $M_t/G/\infty$ models but also serve as approx-

imations (and tools for developing further approximations) for $M_t/G/s/r$ models (with s servers and r extra waiting spaces); see Eick et al. (1993b).

In this paper we continue the investigation of the $M_t/G/\infty$ queue begun in Eick et al. (1993a) by focusing on sinusoidal arrival rate functions. Queues with sinusoidal arrival rates have been considered previously by Jagerman (1975), Rothkopf and Oren (1979), Green and Kolesar (1991) and Green et al. (1991). The special case of sinusoidal arrival rates is especially interesting to understand queues with periodic arrival rates (e.g., daily cycles).

For the $M_t/G/\infty$ queue with sinusoidal arrival rates, we obtain explicit formulas for the mean number of busy servers as a function of time, its extreme values and the time lag between the times when the extremes occur in the arrival rate function and when they occur in the mean number of busy servers. The simple formulas in the special cases of exponential and deterministic service times are especially revealing.

A second purpose of this paper is to study approximations for the time-dependent mean number of busy servers in the $M_t/G/\infty$ model. Since we obtain nice explicit formulas for the case of sinusoidal arrival rates, there obviously is not a great need for approximations here. We are studying approximations in the sinusoidal context to gain insight into how the approximations perform more generally (i.e., for $M_t/G/\infty$ models with general arrival rate functions, for $M_t/G/s/r$ models and for even more general models). We regard the sinusoidal case as a convenient testbed for the approximations.

Here is how this paper is organized. In §2 we define the $M_t/G/\infty$ model and review some of its basic properties; for more discussion and references, see Eick et al. (1993a). In §3 we review some approximations for the mean number of busy servers in an $M_t/G/\infty$ model. In §4 we obtain general results for $M_t/G/\infty$ models with sinusoidal arrival rates. In §§5, 6 and 7 we consider the special cases of exponential, deterministic and hyperexponential service-time distributions, respectively. In §8 we apply the results for sinusoidal arrival rates to treat general periodic arrival rate functions using Fourier series. In §9 we indicate how to calculate the asymptotic sampling variance for the case of a periodic arrival process. Finally, we state some conclusions in §10.

2. The $M_t/G/\infty$ Model

We assume that the $M_t/G/\infty$ model starts empty in the infinite past, which can be formally justified by applying Thorisson (1985). We primarily consider periodic arrival rate functions; then this initialization gives us a dynamic steady state as discussed for $M_t/M/s$ models by Heyman and Whitt (1984).

In general, we assume that λ is nonnegative, measurable and integrable over any bounded interval. For applications, we could also assume that λ is piecewise smooth, i.e., has a continuous derivative everywhere except at finitely many points. Indeed, we use this assumption in §8 for Fourier series.

Let S be a generic service-time random variable and let G be its cumulative distribution function (cdf). Let S_e be a random variable with the associated stationary-excess cdf (or equilibrium-residual-lifetime cdf)

$$G_e(t) \equiv P(S_e \leq t) \equiv \frac{1}{E[S]} \int_0^t G^c(u) du, \quad t \geq 0, \quad (1)$$

where $G^c(t) = 1 - G(t)$; see (16) and (37) of Serfozo (1990). The moments of S_e are related to the moments of S by

$$E[S_e^k] = \frac{E[S^{k+1}]}{(k+1)E[S]}, \quad k \geq 1. \quad (2)$$

Let $Q(t)$ represent the number of busy servers at time t and let $m(t) = E[Q(t)]$. The main $M_t/G/\infty$ result, due to Palm (1943, 1988) and Khintchine (1955, 1960), is that $Q(t)$ has a Poisson distribution with mean

$$\begin{aligned} m(t) &= \int_0^\infty G^c(u) \lambda(t-u) du \\ &= E \left[\int_{t-S}^t \lambda(u) du \right] = E[\lambda(t-S_e)]E[S]. \end{aligned} \quad (3)$$

In Eick et al. (1993a) we review an elegant probabilistic proof of (3) using Poisson random measures, evidently originally due to Prékopa (1958); also see pp. 27–31 of Serfozo (1990). See Eick et al. (1993a) for a review of the literature and Carrillo (1991) for related work in inventory theory.

For interpretation, it is useful to relate the time-dependent mean $m(t)$ in (3) to the instantaneous offered load $\lambda(t)E[S]$, because $m(t)$ equals the instantaneous offered load when the arrival process is homogeneous (by virtue of $L = \lambda W$). For the stationary $M/G/\infty$ model, the steady-state mean $m(\infty)$ thus depends on the service-time distribution only through its mean; i.e., the $M/G/\infty$ model has the familiar insensitivity property. From (3) it is clear that this is not true for the $M_t/G/\infty$ model; the entire service-time distribution plays a role.

The last expression for $m(t)$ in (3) says that to obtain $m(t)$ we replace $\lambda(t)$ in the instantaneous offered load by an appropriate weighted average of λ before time t . Since we average before time t , we see that the congestion as described by $m(t)$ lags behind $\lambda(t)$; e.g., if λ is unimodal with a unique maximum, then the peak of m will come after the peak of λ ; see Corollary 2.5 of Eick et al. (1993a).

The second expression for $m(t)$ in (3) also has a nice interpretation. The integral describes the cumulative arrival rate during a service period before time t . The expectation is then the expected number of arrivals during a random service period before t . The probabilistic proof shows that this interpretation is justified.

It is significant that the operator mapping the function λ into the function m in (3) is a *linear operator*: see Theorem 2.10 of Eick et al. (1993a). Thus, we can regard the congestion measure m as a response of a *linear system* (the queueing model) to an input signal (the arrival rate function λ); e.g., see Chapter 2 of Ziemer and Tranter (1976). This linearity occurs in large part because different customers do not interfere with each other. The linear structure in the $M_t/G/\infty$ model explains why the $M_t/G/\infty$ model is much easier to analyze than other $M_t/G/s/r$ queues.

It is not necessary to be familiar with linear system theory to read this paper, but linear system theory adds additional insight. The first expression for $m(t)$ in (3) is the standard *superposition integral* in linear system theory; see p. 53 of Ziemer and Tranter (1976). The *impulse response function* (i.e., the response to the system to an impulse applied at time $t = 0$) is G^c . The superposition integral shows that the response m is the convolution of the impulse response function G^c with the signal λ . By (1), the impulse function is just $E[S]$ times the probability density of S_e . Linear system theory provides a nice framework for interpreting the results in this paper. An important role is played by the *transfer function*, which is the Fourier transform of the impulse response function; see (11) below.

3. Approximations for the Mean Function

Since the distribution of $Q(t)$ is Poisson for each t , to describe this distribution as a function of time it suffices to focus on the mean function m in (3). For any arrival rate function λ and any service-time cdf G , $m(t)$ is easily calculated from (3), using numerical integration if necessary. However, to *understand* the behavior of $M_t/G/\infty$ systems, and perhaps to calculate more quickly, it is useful to consider various approximations for $m(t)$. We review some basic approximations for $m(t)$ here; they are summarized in Table 1.

An obvious approximation strategy, commonly applied in practice, is to act as if the arrival process were homogeneous. The *simple stationary approximation* (SSA) uses the stationary $M/G/\infty$ formula with a long-run average arrival rate $\bar{\lambda}$, assuming that the long-run average is well defined, as it is in the periodic case. An alternative scheme is the *pointwise stationary approxi-*

Table 1 **A Summary of Basic Approximations for the Mean Function m**

SSA	$\bar{\lambda}E[S]$
PSA	$\lambda(t)E[S]$
LIN-S	$\lambda(t - E[S_e])E[S]$
LIN-D	$\lambda(t)E[S] - \frac{\lambda^{(1)}(t)E[S^2]}{2}$
QUAD-S	$\lambda(t - E[S_e])E[S] + \frac{\lambda^{(2)}(t)}{2} \text{Var}(S_e)E[S]$
QUAD-D	$\lambda(t)E[S] - \frac{\lambda^{(1)}(t)E[S^2]}{2} + \frac{\lambda^{(2)}(t)E[S^3]}{6}$
CUBIC-D	$\text{QUAD-D} - \frac{\lambda^{(3)}(t)E[S^4]}{24}$

mation (PSA), which calculates $m(t)$ as if the arrival process were homogeneous with the instantaneous arrival rate $\lambda(t)$. Thus, for $m(t)$, SSA is $\bar{\lambda}E[S]$, while PSA is $\lambda(t)E[S]$. (For finite-capacity systems, these stationary approximations could be unstable, but that cannot occur here.) Note that PSA is a function of time, whereas SSA is a constant. The corresponding approximations for multi-server delay systems are examined by Green and Kolesar (1991) and Green et al. (1991).

In Eick et al. (1993a) various *polynomial approximations* for $m(t)$ were introduced, which are based on assuming that the arrival rate function λ is a polynomial or can be approximated by a polynomial. The approximations LIN-S (linear with time shift) and QUAD-S (quadratic with time shift) are based on the formula

$$m(t) = \lambda(t - E[S_e])E[S] + \frac{\lambda^{(2)}(t)}{2} \text{Var}(S_e)E[S], \quad (4)$$

where $\lambda^{(k)}(t)$ is the k th derivative of λ at t . Formula (4) is valid if λ is quadratic; see (7) and (15) of Eick et al. (1993a). The approximations LIN-D, QUAD-D and CUBIC-D (D for derivatives) are based on Taylor series expansions for λ ; i.e.,

$$m(t) = \sum_{j=0}^n (-1)^j \frac{\lambda^{(j)}(t)E[S^{j+1}]}{(j+1)!} + R_n(t), \quad (5)$$

where $R_n(t)$ is a remainder term; see (8), (16) and Theorem 3.2 of Eick et al. (1993a). These approximations all appear in Table 1.

As indicated above, we will examine how these approximations perform in the special case of an $M_t/G/\infty$ model with sinusoidal arrival rate function. We obtain a very clear picture because we can obtain convenient explicit expressions for $m(t)$ as well as the approximations.

4. General Results for Sinusoidal Arrival Rates

Now we consider the sinusoidal arrival rate function

$$\lambda(t) \equiv \bar{\lambda} + \beta \sin(\gamma t) \equiv \bar{\lambda} + \bar{\lambda}\alpha \sin(2\pi t/\psi), \quad (6)$$

for positive constants $\bar{\lambda}$, $\beta \equiv \bar{\lambda}\alpha$ and $\gamma \equiv 2\pi/\psi$, where $0 < \alpha < 1$. The second representation in (6) is convenient for interpretation, because $\bar{\lambda}$ is the *average arrival rate*, α is the *relative amplitude* and ψ is the *cycle length or period*. We call γ the *frequency*. In the context of (6), the arrival process is characterized by *three parameters*, e.g., the triple $(\bar{\lambda}, \beta, \gamma)$, the triple $(\bar{\lambda}, \alpha, \psi)$ or the triple $(\bar{\lambda}, \alpha, \gamma)$. Of course, these parameters should be interpreted relative to the mean service time $E[S]$. There is one degree of freedom for choosing the measuring units. If we choose measuring units so that $E[S] = 1$, we speak of $\bar{\lambda}$ as the *relative average arrival rate*, ψ as the *relative cycle length* and γ as the *relative frequency*.

As pointed out by Green et al., a key role is played by the *relative cycle length*. In many applications, a cycle represents a day. Table 2 displays values of the relative cycle length and relative frequency as a function of the mean service time, assuming a daily cycle. For example, if we think of a daily cycle applying to telephone calls with a mean holding time of five minutes, then $\psi = 288$ and $\gamma = 0.022$. More detailed modeling of telephone traffic might lead to two cycles over the twelve hour period from 8 am to 8 pm. Then we have $\psi = 72$ and $\gamma = 0.087$. In this context, Palm (1943, 1988) suggested that PSA should be a pretty good approximation, and this will be substantiated for the mean $m(t)$ in the $M_t/G/\infty$ model here. More generally, we will investigate the error in PSA as a function of the parameters $\bar{\lambda}$, α , γ and the service-time cdf G , assuming that $ES = 1$.

We begin by obtaining an explicit expression for m .

THEOREM 4.1. For sinusoidal λ as in (6),

$$m(t) = \bar{\lambda}E[S] + \beta(\sin(\gamma t)E[\cos(\gamma S_e)] - \cos(\gamma t)E[\sin(\gamma S_e)])E[S]. \quad (7)$$

Table 2 The Relative Cycle Length ψ and the Relative Frequency γ as a Function of Mean Service Time for a Daily Cycle.

(The Relative Cycle Length and Relative Frequency Are the Cycle Length and Frequency Computed With Measuring Units So That $E[S] = 1$.)

Mean Service Time $E[S]$	Relative Cycle Length ψ	Relative Frequency γ
1 second	86,400	7.27×10^{-5}
1 minute	1,440	4.36×10^{-3}
5 minutes	288	0.0218
10 minutes	144	0.0436
30 minutes	48	0.131
1 hour	24	0.262
4 hours	6	1.05
12 hours	2	3.14
1 day	1	6.28
7 days	1/7	44.0

PROOF. From (3), we immediately obtain

$$m(t) = (\bar{\lambda} + \beta E[\sin(t - S_e)])E[S].$$

The conclusion then follows from the sine addition formula

$$\sin(x - y) = \sin x \cos y - \cos x \sin y. \quad \square \quad (8)$$

REMARK (4.1). An easy alternate proof, which can provide additional insight, is via Fourier transforms; see pp. 58–59 of Ziemer and Tranter (1976). Corollary 4.2 below is also familiar in that context. \square

From (7), we see that m is periodic with period (cycle length) $\psi = 2\pi/\gamma$ just like λ . Moreover, the long-run average (and average over one cycle) is

$$\bar{m} = \lim_{t \rightarrow \infty} t^{-1} \int_0^t m(s) ds = \bar{\lambda}E[S]. \quad (9)$$

Formula (9) implies that the approximations SSA and PSA are both *exact in an average sense*. Note that this is in distinct contrast with the behavior of queues with finitely many servers. Experience related to Ross's (1978) conjecture (e.g., Rolski (1989) and Green et al. (1991)) indicates that with finitely many servers the SSA value *underestimates* congestion. Formula (9) suggests that SSA should perform better in this average sense as the number of servers increases and the model becomes more like an infinite-server model.

Of course, SSA completely fails to say anything about the time-dependent behavior of m . For this, we may

turn to PSA. However, because of the averaging (3), the extreme values of m will necessarily be relatively less extreme than the extreme values $\bar{\lambda} \pm \beta$ of λ . Hence, while SSA overestimates the true averaging effect at a particular time, PSA underestimates it.

Note that the extreme values of λ occur at times

$$t_\lambda = \frac{\pi}{2\gamma} + \frac{\pi n}{\gamma} \quad (10)$$

for integer n , with the maximum occurring when n is even. We now apply (7) and elementary calculus to describe the locations and values of the extremes of m . Here we will use elementary properties of complex numbers. Recall that if $z = x + iy$, where $i = \sqrt{-1}$, then the *modulus* of z is $|z| = (x^2 + y^2)^{1/2}$ and the *argument* of z , denoted by $\arg(z)$, is the angle θ between the x -axis and the vector (x, y) in the plane, i.e., $\tan \theta = y/x$. Recall that $e^z = e^x(\cos y + i \sin y)$.

COROLLARY 4.2. *The extreme values of m occur at times*

$$\begin{aligned} t_m &= \gamma^{-1} \tan^{-1} (-E[\cos \gamma S_e]/E[\sin \gamma S_e]) \\ &= t_\lambda + \gamma^{-1} \tan^{-1} (E[\sin \gamma S_e]/E[\cos \gamma S_e]) \\ &= t_\lambda + \gamma^{-1} \arg (E[e^{i\gamma S_e}]), \end{aligned} \quad (11)$$

and the extreme values are

$$\begin{aligned} m(t_m) &= \bar{\lambda}E[S] \pm \beta |E[e^{i\gamma S_e}]| E[S] \\ &= \bar{\lambda}E[S] \pm \beta (E[\cos \gamma S_e])^2 \\ &\quad + (E[\sin \gamma S_e])^2)^{1/2} E[S], \end{aligned} \quad (12)$$

so that

$$|m(t_m) - \bar{\lambda}E[S]| \leq \beta E[S].$$

REMARK (4.2). In linear system theory we obtain Corollary 4.2 simply by noting that the amplitude response function and the phase shift function are the modulus and argument, respectively, of the transfer function at the single relevant frequency γ ; see pp. 33 and 55 of Ziemer and Tranter (1976). \square

From (11), we see that the lag $(t_m - t_\lambda)$ depends on λ in (6) only via the frequency γ . It is evident that the arrival rate function λ is symmetric about its extremes, i.e., $\lambda(t_\lambda + t) = \lambda(t_\lambda - t)$ for all t . It is interesting that the mean function m has this symmetry property too.

THEOREM 4.3. *$m(t_0 + t) = m(t_0 - t)$ for all t if and only if $t_0 = t_m$ for t_m in (11).*

PROOF. Apply the sine addition formula (8) with (7) to calculate $m(t_0 + t) - m(t_0 - t)$. Then observe that the value is zero if and only if $t_0 = t_m$. \square

Theorem 4.3 implies that congestion (as measured by m) decreases after its peak just as quickly as it increases before the peak. As shown by Green et al. (1991), with finitely many servers, congestion tends to increase faster before its peak than it decreases afterwards. Theorem 4.3 indicates that this effect disappears as the number of servers increases.

It is instructive to see what happens as the frequency γ gets very small or very large. In the following we index λ and m by γ . Below we establish the asymptotic validity of PSA when γ gets small and the asymptotic validity of SSA when γ gets large. These results are special cases of what can be established more generally. For example, in the case of exponential service times Theorem 4.3 is a corollary to Theorem 1 of Whitt (1991), but the direct proof here is much easier.

Intuitively, we should anticipate that PSA should perform well as $\gamma \rightarrow 0$, because then the cycles are very long, i.e., λ changes slowly relative to the mean service time. However, the desired asymptotic behavior as $\gamma \rightarrow 0$ requires some care to state. For any fixed t , $\lambda_\gamma(t) \rightarrow \lambda(0) = \bar{\lambda}$ and $m_\gamma(t) \rightarrow \bar{\lambda}E[S]$, so that SSA is also asymptotically correct in this sense. However, we do not really want to consider a fixed time t , but instead a fixed position within a cycle. To achieve this, we must scale time in m_γ . In the system with frequency γ we want to consider times t/γ . This can also be achieved by having uniform convergence in t as $\gamma \rightarrow 0$.

THEOREM 4.4. *If $\gamma \rightarrow 0$, then $m_\gamma(t) - \lambda_\gamma(t)E[S] \rightarrow 0$ and $m_\gamma(t/\gamma) \rightarrow \lambda_1(t)E[S]$ uniformly in t .*

PROOF. Since $\cos \gamma x \rightarrow 1$ and $\sin \gamma x \rightarrow 0$ as $\gamma \rightarrow 0$ for all x , the results follow easily from Theorem 4.1. In particular, $E[\cos(\gamma S_e)] \rightarrow 1$ and $E[\sin(\gamma S_e)] \rightarrow 0$ by the bounded convergence theorem. \square

In contrast to the situation when $\gamma \rightarrow 0$, we should anticipate that PSA will perform poorly as $\gamma \rightarrow \infty$. For example, the number of arrivals in any interval $[a, b]$ is Poisson with mean $\int_a^b \lambda(t)dt$, so that the crucial quantity is the *integral* of λ , not λ itself. (Also see (3) and the proof below.) If λ is suitably smooth, then

$$\int_a^b \lambda(t) dt \approx \lambda(a)(b - a),$$

but if λ oscillates rapidly, then any specific values of $\lambda(t)$ will be misleading. In contrast, SSA is asymptotically correct as $\gamma \rightarrow \infty$. As $\gamma \rightarrow \infty$, there will be many complete cycles during a service time, so that the relevant arrival rate is the long-run average rate $\bar{\lambda}$.

THEOREM 4.5. *If $\gamma \rightarrow \infty$, then $m_\gamma(t) \rightarrow \bar{\lambda} E[S]$ uniformly in t .*

PROOF. The arrival process with rate function λ_γ converges in distribution to a homogeneous Poisson process with constant rate $\bar{\lambda}$ as $\gamma \rightarrow \infty$, because the compensators converge, i.e.,

$$\begin{aligned} & \int_0^t [\bar{\lambda} + \beta \sin(\gamma u)] du \\ &= \bar{\lambda}t + \frac{\beta}{\gamma} \cos(\gamma t) \rightarrow \bar{\lambda}t \quad \text{as } \gamma \rightarrow \infty; \end{aligned}$$

see Theorem 5.7 of Serfozo (1990). Then a continuity argument shows that $Q_\gamma(t)$ converges in distribution to $Q(t)$ for each t as $\gamma \rightarrow \infty$, e.g., by a minor modification of Whitt (1974). Finally, $m_\gamma(t) \rightarrow m(t)$ as $\gamma \rightarrow \infty$ since the distributions are all Poisson. As a consequence (or by a direct argument), $E[\cos(\gamma S_e)] \rightarrow 0$ and $E[\sin(\gamma S_e)] \rightarrow 0$ as $\gamma \rightarrow 0$ in (8), but then $m_\gamma(t) \rightarrow \bar{\lambda}E[S]$ as $\gamma \rightarrow \infty$ uniformly in t . \square

Of course, Theorems 4.4 and 4.5 do not tell when γ is sufficiently small or large for the limits to be realized as good approximations in practice. However, we can write down explicit expressions for the errors:

$$\begin{aligned} e_{\text{PSA}}(t) &\equiv \lambda(t)E[S] - m(t) \\ &= \bar{\lambda}E[S](\sin(\gamma t)(1 - E[\cos(\gamma S_e)]) \\ &\quad + \cos(\gamma t)E[\sin(\gamma S_e)]) \quad \text{and} \quad (13) \end{aligned}$$

$$\begin{aligned} e_{\text{SSA}}(t) &\equiv \bar{\lambda}E[S] - m(t) \\ &= \bar{\lambda}E[S](\sin(\gamma t)E[\cos(\gamma S_e)] \\ &\quad + \cos(\gamma t)E[\sin(\gamma S_e)]). \quad (14) \end{aligned}$$

From (13) and (14), we see that both errors $e_{\text{PSA}}(t)$ and $e_{\text{SSA}}(t)$ are proportional to the average arrival rate $\bar{\lambda}$. Consequently, the relative errors $e_{\text{PSA}}(t)/m(t)$ and $e_{\text{SSA}}(t)/m(t)$ depend on $(\bar{\lambda}, \alpha, \psi)$ only through (α, ψ) .

5. Exponential Service Times

Consider the special case of an exponential service-time distribution with mean 1. Since then $E[\sin(\gamma S_e)] = \gamma/(1 + \gamma^2)$ and $E[\cos(\gamma S_e)] = 1/(1 + \gamma^2)$. From (8) we see that

$$m(t) = \bar{\lambda} + \frac{\beta}{1 + \gamma^2} [\sin \gamma t - \gamma \cos \gamma t]. \quad (15)$$

Formula (15) can also be verified by applying Corollary 2.8 of Eick et al. (1993a) which says that for any $M_t/M/\infty$ model

$$m'(t)E[S] = \lambda(t)E[S] - m(t).$$

This differential equation in turn is just an expression for a more general conservation law stating that $m'(t)$ is the arrival rate minus the departure rate; see Theorem 2.6 of Eick et al. (1993a). For the Markovian $M_t/M/\infty$ model, it is obvious that the conservation law takes the form $m'(t) = \lambda(t) - \mu m(t)$, where $\mu = 1/E[S]$. The differential equation implies that the PSA approximation $\lambda(t)E[S]$ coincides with $m(t)$ precisely when $m'(t) = 0$.

REMARK (5.1). From the perspective of linear system theory, this differential equation shows that the $M_t/M/\infty$ model corresponds to the low-pass RC filter; see pp. 50 and 55 of Ziemer and Tranter (1976). \square

From (11), and the fact that we have fixed the measuring units by setting $E[S] = 1$, we see that $m'(t) = 0$ at times

$$t_m = t_\lambda + \cot^{-1}(1/\gamma)/\gamma, \quad (16)$$

i.e., the time lag in the peak of m after the peak of λ is $L(\gamma) \equiv \cot^{-1}(1/\gamma)/\gamma$ and the phase shift (relative lag in the peak per cycle) is

$$\phi(\gamma) \equiv \frac{L(\gamma)}{(2\pi/\gamma)} = \frac{\cot^{-1}(1/\gamma)}{2\pi}. \quad (17)$$

The time lag $L(\gamma)$ is decreasing in γ , going from $L(0) = 1$ to $L(\infty) = 0$, while the phase shift in (17) is increasing in γ , going from $\phi(0) = 0$ to $\phi(\infty) = \frac{1}{4}$.

From (12), we see that extreme values of m in this case are

$$m(t_m) = \bar{\lambda} \pm \frac{\beta}{\sqrt{1 + \gamma^2}}. \quad (18)$$

From (18) we see that the extremes are decreasing in γ , with the limits as $\gamma \rightarrow 0$ and $\gamma \rightarrow \infty$ being consistent with Theorems 4.4 and 4.5. Moreover, we have a convenient quantitative *tight bound* for SSA, which is the approximation provided by Theorem 4.5:

$$|e_{SSA(t)}| \equiv |m(t) - \bar{\lambda}| \leq \beta / \sqrt{1 + \gamma^2} \quad \text{for all } t. \quad (19)$$

The relative error is then

$$\frac{e_{SSA(t)}}{m(t)} \approx \frac{e_{SSA(t)}}{\bar{\lambda}E[S]} \approx \frac{\alpha}{\sqrt{1 + \gamma^2}} \quad (20)$$

for α and γ in (6). Assuming that the relative amplitude α is not small, we see from (20) that the maximum relative error in SSA as a pointwise approximation is small only if $\gamma \gg 1$ (γ is large compared to 1), which only occurs if the mean service time is greater than one day in the setting of Table 2. Thus, we conclude that SSA is typically a poor pointwise approximation for periodic arrival rates.

EXAMPLE 5.1. To illustrate, consider an $M_t/M/\infty$ model with exponential service times having mean 1 and arrival rate function $\lambda(t) = 10 + 5 \sin(t)$, i.e., $\bar{\lambda} = 10$, $\alpha = 0.5$ and $\gamma = 1$. Figure 1 displays λ , which coincides with PSA, and m over a time interval of four cycles (25.13). As indicated above, $m'(t) = 0$ precisely where $\lambda(t) = m(t)$. From (18), the relative amplitude of m is

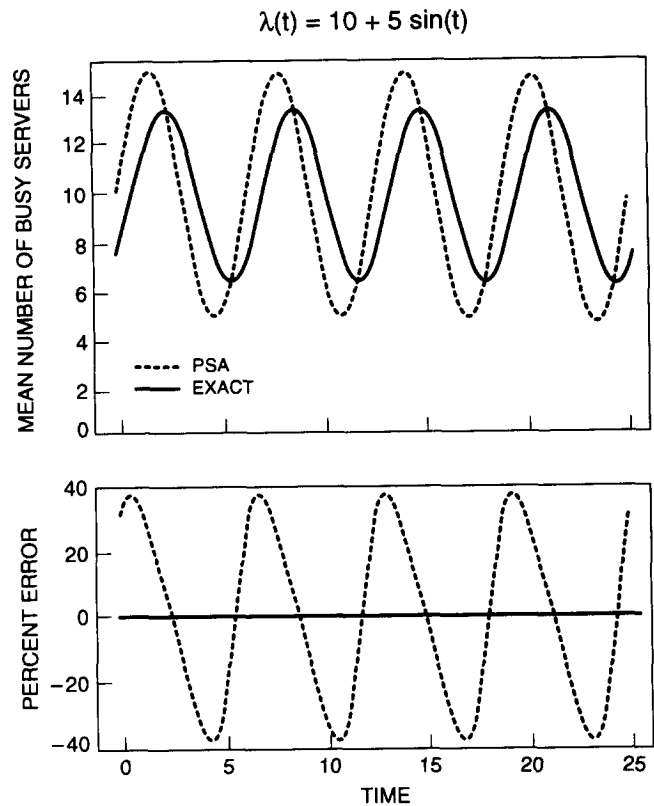
$$a / \sqrt{1 + \gamma^2} = \sqrt{2} / 4 \approx 0.353$$

as opposed to $\alpha = 0.500$ for λ . From (16), the time lag in the extremes is $\cot^{-1}(1) = \pi/4 \approx 0.785$.

The percent error in PSA in (13) is plotted in Figure 1 as well as the mean number of busy servers. Note that the maximum percent error tends to be somewhat larger than we might suspect from only looking at the plots of λ and m . The maximum error occurs less than halfway between the extreme values of m (where there is no error). Also note that the maximum error occurs before the extreme in λ . \square

Note that λ in (6) has bounded derivatives of all orders, so that we can apply Theorem 3.2 of Eick et al. (1993a) to obtain Taylor series expansions for $m(t)$. In this case, the expansion is absolutely convergent if and only if $\gamma < 1$. Moreover,

Figure 1 The Comparison of $\lambda(t) = \text{PSA}$ with $m(t)$ in Example 5.1



$$m_{2n+1}(t) = \bar{\lambda} \left(1 + \alpha \left(\sum_{j=0}^n (-1)^j (\gamma^2)^j \right) \times (\sin \gamma t - \gamma \cos \gamma t) \right), \quad n \geq 0, \quad (21)$$

which is of the same form as (13) except that the term

$$(1 + \gamma^2)^{-1} = \sum_{j=0}^{\infty} (-1)^j (\gamma^2)^j$$

in (13) is approximated by the finite sum $\sum_{j=0}^n (-1)^j \times (\gamma^2)^j$. In this case with S exponential and $E[S] = 1$, the remainder term is simply

$$R_n(t) = E[\lambda^{(n+1)}(t - S)], \quad (22)$$

because $S_e = S$ and $E[S^n] = n!$. From (22), we see that

$$\sup_{t \geq 0} |R_n(t)| = \bar{\lambda} \alpha \gamma^{n+1}. \quad (23)$$

The first four terms of the Taylor series expansion are PSA, LIN-D, QUAD-D and CUBIC-D in Table 1. In

particular, LIN-D and CUBIC-D are (21) for $n = 0$ and 1, respectively.

In this case the relative error in PSA is

$$\frac{e_{PSA}(t)}{m(t)} = \frac{\gamma^2 \sin(\gamma t) + \gamma \cos \gamma t}{1 + \gamma^2 + \alpha[\sin \gamma t - \gamma \cos \gamma t]} \quad (24)$$

Figure 2 displays contours of maximum (over t) percent PSA error as a function of the two parameters γ and α . For example, corresponding to $\gamma = 0.087$ for the telephone example mentioned in the beginning of §4, the maximum percent error in $m(t)$ using PSA increases from about 1% to 10% as α increases from 0.2 to 0.8. A realistic value of α might be 0.4, which corresponds to a maximum percent error in $m(t)$ from PSA of 4%.

EXAMPLE 5.2. We now modify Example 1 by considering smaller and more realistic relative frequency factors γ . In particular, we let $\bar{\lambda} = 10$ and $\alpha = 0.5$ as before, but now we let $\gamma = 0.5, 0.2$ and 0.02 . These values correspond to mean service times of about 115 minutes, 46 minutes and 5 minutes, respectively, in a daily cycle. Figures 3–5 compare several approximations with the exact values of $m(t)$ for this example.

Figure 2 The Contours of Maximum Percent Error in the PSA Approximation for $m(t)$ as a Function of the Relative Amplitude α and the Relative Frequency γ when the Service Time is Exponential with Mean 1

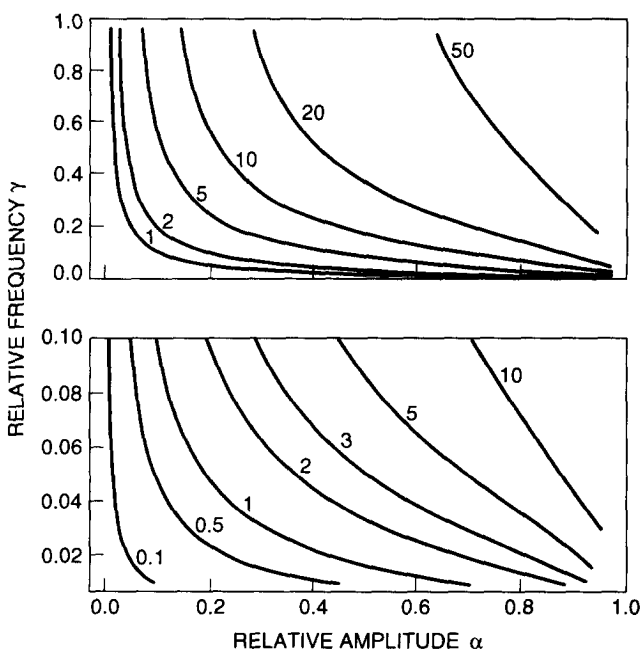
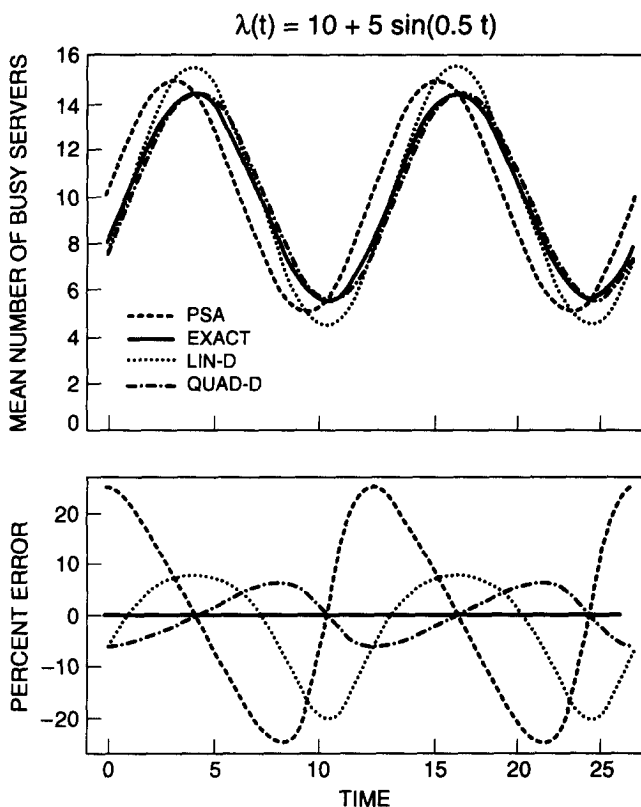


Figure 3 A Comparison of Approximations with Exact Values when $\gamma = 0.5$ in Example 5.2



As in Figure 1, in Figure 3 for $\gamma = 0.5$ we see the possible difficulties in interpretation. First, the periodic plots of the means tend to hide the actual errors. For example, from the periodic plots of the means, we might not realize that the maximum percent error in PSA is as much as 23% when $\gamma = 0.5$. On the other hand, the errors may exaggerate the true differences in the mean functions, because the errors measure vertical distance only. The mean functions could be considered closer if we used a different metric.

For $\gamma = 0.2$ and 0.02 , we only give the percent errors, because the curves are too close in direct plots. We plot the percent errors with and without PSA. We see that the error in PSA is much greater than the other errors and that the error in LIN-D is much greater than QUAD-D.

EXAMPLE 5.1 (REVISITED). We have noted that the Taylor series for $m(t)$ converges absolutely if and only if $\gamma < 1$. Thus, we should not expect the polynomial

Figure 4 A Comparison of Approximations with Exact Values when $\gamma = 0.2$ in Example 5.2

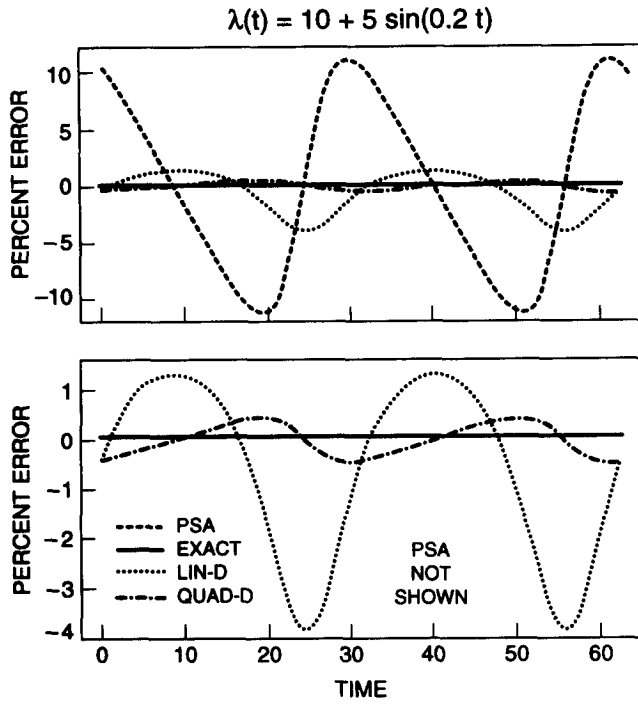
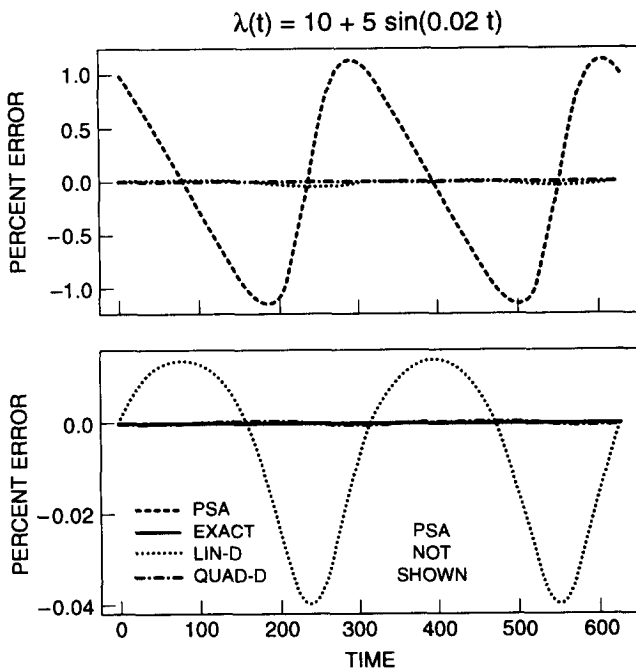


Figure 5 A Comparison of Approximations with Exact Values when $\gamma = 0.02$ in Example 5.2



approximations to perform well when $\gamma \geq 1$. This is illustrated by Figure 6, which compares approximations to the exact values of $m(t)$ when $\bar{\lambda} = 10$, $\alpha = 0.5$ and $\gamma = 1$, as in Example 5.1.

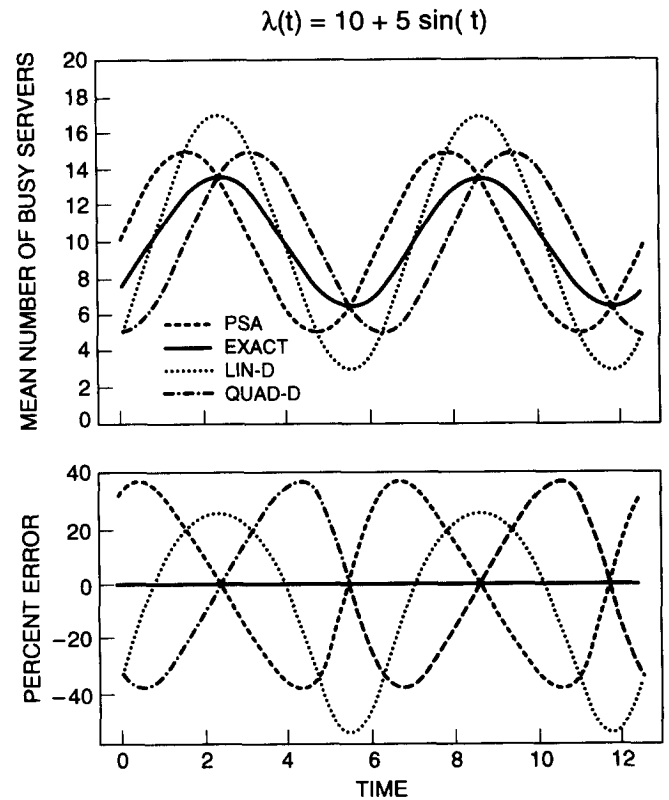
6. Deterministic Service Times

Now consider the case of deterministic service times, all assuming the value 1. As indicated before, with the stationary $M/G/\infty$ model, the steady-state mean $m(\infty)$ is $\bar{\lambda}E[S]$, independent of the service-time distribution (which is quite different from the effect of the service-time distribution in the stationary $M/G/1$ model). However, the service-time distribution beyond its mean affects the mean function m in the $M_t/G/\infty$ model.

From Theorem 4.1, we see that

$$m(t) = \bar{\lambda} + \frac{\beta}{\gamma} [\cos(\gamma(t-1)) - \cos(\gamma t)]. \quad (25)$$

Figure 6 A Comparison of Approximations with Exact Values when $\gamma = 1.0$ in Example 5.1 Revisited



Formula (25) can also be verified by applying Corollary 2.7 of Eick et al. (1993a), which states that $m'(t) = \lambda(t) - \lambda(t - E[S])$ for any $M_t/D/\infty$ model. This corollary also shows that $m'(t) = 0$ at times

$$t_m = t_\lambda + \frac{1}{2} \quad \text{and} \quad (26)$$

$$m(t_m) = \bar{\lambda} \pm \frac{2\beta}{\gamma} \sin\left(\frac{\gamma}{2}\right). \quad (27)$$

Since $(\sin x)/x \leq 1$ for all x , $|m(t_m)| \leq \bar{\lambda} + \beta$, as required. From (26), we see that in contrast to the case of exponential service times, in which the lag $L(\gamma)$ decreases from 1 to 0 and the phase shift $\phi(\gamma)$ increases from 0 to $\frac{1}{4}$ as γ increases, here there is a fixed lag $L(\gamma) = \frac{1}{2}$ and the phase shift is $\phi(\gamma) = \gamma/4\pi$, $0 < \gamma < 4\pi$, which increases from 0 to 1 as γ increases from 0 to 4π .

7. Hyperexponential Service Times

Now suppose that the distribution of S is hyperexponential H_k (a mixture of k exponentials), i.e.,

$$G^c(t) = \sum_{i=1}^k p_i e^{-\mu_i t}, \quad t \geq 0, \quad (28)$$

where $E[S] = \sum_{i=1}^k (p_i/\mu_i) = 1$. Then, from Theorem 2.10 of Eick et al. (1993a) and (15), we obtain

$$m(t) = \bar{\lambda} + \beta \sum_{i=1}^k \frac{p_i}{\mu_i} \left[\frac{\gamma \mu_i}{(\mu_i + \gamma)^2} \sin(\gamma t) - \frac{\mu_i^2}{(\mu_i + \gamma)^2} \cos(\gamma t) \right]. \quad (29)$$

EXAMPLE 7.1. Consider an H_2 distribution with balanced means ($p_1/\mu_1 = p_2/\mu_2 = 0.5$) and squared coefficient of variation (variance divided by the square of the mean) $c_s^2 = 5$. Then $p_1 = 0.908$, $\mu_1 = 1.816$, $p_2 = 0.092$ and $\mu_2 = 0.184$; see (3.7) of Whitt (1982). If $\gamma = 1$, then $E[\sin(\gamma S_e)] = 0.126$ and $E[\cos(\gamma S_e)] = 0.220$, so that

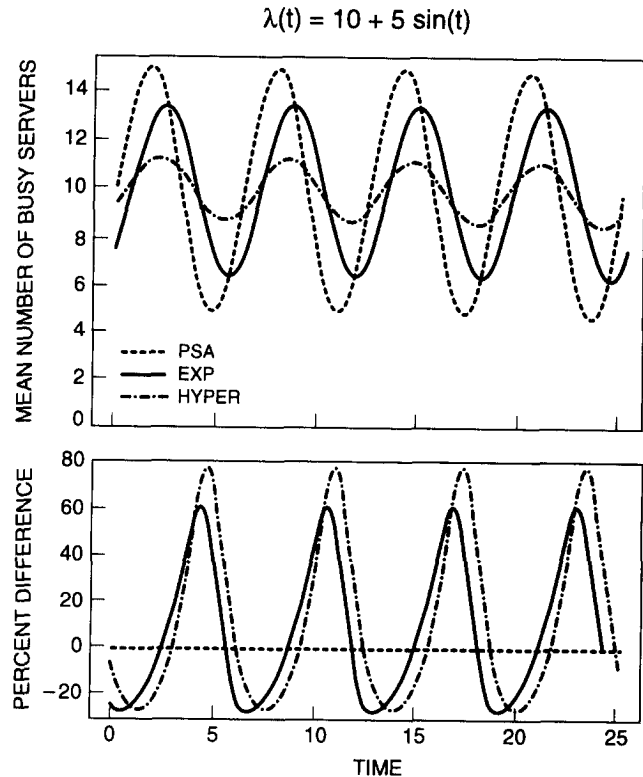
$$m(t) = \bar{\lambda} + \beta(0.220 \sin t - 0.126 \cos t).$$

In contrast, for an exponential service time with mean 1,

$$m(t) = \bar{\lambda} + \beta(0.500 \sin t - 0.500 \cos t).$$

These mean functions are compared in Figure 7.

Figure 7 A Comparison of $m(t)$ for Exponential and Hyperexponential Service-time Distributions in Example 7.1



8. General Periodic Arrival Rate Functions

In this section we assume that λ is a general periodic function on $[0, 2\pi/\gamma]$, and analyze the associated $M_t/G/\infty$ model by applying Fourier series together with the results in §4. Note that this periodic case essentially covers a general arrival rate function on a finite interval, because any such arrival rate function can be extended to a periodic function. The only difficulty is the end effect at the left boundary in the aperiodic function, which can usually be represented by appropriately modifying the periodic function. In particular, if the given arrival rate function on the interval $[b, c]$ is initialized at b as if there were a warmup period during $[a, b]$, then we construct our periodic function on $[a, c]$ and use the results from $[b, c]$.

Since λ is periodic on $[0, 2\pi/\gamma]$, it follows from (3) that m is also periodic on $[0, 2\pi/\gamma]$. To treat general λ , we assume that λ can be approximated by the partial sums of its Fourier series, i.e.,

$$\lambda_n(t) = a_0 + \sum_{k=1}^n (a_k \sin k\gamma t + b_k \cos k\gamma t), \quad (30)$$

where

$$a_k = \frac{1}{\pi} \int_0^{2\pi} \lambda(t) \sin k\gamma t dt \quad \text{and} \\ b_k = \frac{1}{\pi} \int_0^{2\pi} \lambda(t) \cos k\gamma t dt. \quad (31)$$

To guarantee convergence of λ_n in (30) as $n \rightarrow \infty$, we assume that λ is piecewise smooth on $[0, 2\pi/\gamma]$, i.e., λ has a continuous derivative except at finitely many points where λ and its derivative may have simple jump discontinuities. Then $\lambda_n(t) \rightarrow \lambda(t)$ as $n \rightarrow \infty$ for each t that is a point of continuity of λ and

$$\lambda_n(t) \rightarrow [\lambda(t+) + \lambda(t-)]/2$$

at each point of discontinuity. Moreover, the convergence is uniform if λ is continuous everywhere; see pp. 19, 81 of Tolstov (1976). From (3), we see that the mean function m_n associated with λ_n converges uniformly to m when λ_n converges uniformly to λ ; see Theorem 2.10 of Eick et al. (1993a).

Moreover,

$$m_n(t) = a_0 + \sum_{k=1}^n a_k m_{k1}(t) + b_k m_{k2}(t), \quad (32)$$

where $m_{k2}(t) = m_{k1}(t + \pi/2\gamma)$ because $\cos(kt) = \sin(kt + \pi/2\gamma)$ and $m_{k1}(t)$ is the formula from Theorem 4.1 with γ replaced by $k\gamma$.

EXAMPLE 8.1. If the service times are exponential with mean 1, then we can combine (15) and (32) to obtain

$$m_n(t) = a_0 + \sum_{k=1}^n \frac{a_k}{1 + (k\gamma)^2} \\ \times [\sin(k\gamma t) - k\gamma \cos(k\gamma t)] \\ + \sum_{k=1}^n \frac{b_k}{1 + (k\gamma)^2} [\sin(k\gamma t + \pi/2) \\ - k\gamma \cos(k\gamma t + \pi/2)] \\ = a_0 + \sum_{k=1}^n \frac{(a_k + kb_k\gamma)}{1 + (k\gamma)^2} \sin(k\gamma t) \\ + \sum_{k=1}^n \frac{(b_k - ka_k\gamma)}{1 + (k\gamma)^2} \cos(k\gamma t). \quad (33)$$

From (33), we see that the terms with high k have relatively little influence because of the term $(1 + (k\gamma)^2)$ in the denominator of the coefficients. Thus, m might be well approximated by fewer terms from its Fourier series than λ .

9. The Asymptotic Sampling Variance

In this section we indicate how to calculate the asymptotic sampling variance for the case of a periodic arrival rate. This is useful for determining confidence intervals for the long-run average. For this purpose, we establish a preliminary lemma. We use the iterate of the stationary-excess operator in (1), i.e., $S_e^{(k)} = (S_e^{(k-1)})_e$.

LEMMA 9.1. For an $M_t/G/\infty$ system,

$$2 \int_0^\infty \text{Cov}[Q(t), Q(t+u)] du = E[\lambda(t - S_e^{(2)})]E[S^2].$$

PROOF. From Theorem 1.2 of Eick et al. (1991), we have

$$\int_0^\infty \text{Cov}[Q(t), Q(t+u)] du \\ = \int_0^\infty E \left[\int_{t-(s-s)^-}^t \lambda(s) ds \right] du \\ = \int_0^\infty \int_{-\infty}^t \lambda(s) P(S > t+u-s) ds du \\ = \int_{-\infty}^t \lambda(s) \int_0^\infty P(S > t+u-s) du ds \\ = E[S] \int_{-\infty}^t \lambda(s) P(S_e > t-s) ds \\ = E[S]E[S_e]E[\lambda(t - S_e^{(2)})] \\ = E[\lambda(t - S_e^{(2)})]E[S^2]/2. \quad \square$$

THEOREM 9.1. If $\lambda(t)$ is a periodic function with period T , then the asymptotic sampling variance of $Q(t)$ is

$$\lim_{t \rightarrow \infty} t^{-1} \text{Var} \left[\int_0^t Q(s) ds \right] = \frac{1}{T} \int_0^T E[\lambda(t - S_e^{(2)})] dt E[S^2].$$

PROOF. If $Q(t)$ were a stationary process, then we would have

$$\lim_{t \rightarrow \infty} t^{-1} \text{Var} \left[\int_0^t Q(s) ds \right] = 2 \int_0^\infty \text{Cov}[Q(0), Q(t)] dt,$$

but it is not. However, we can make $\{Q(t) : t \geq 0\}$ stationary by randomizing the place in the cycle at time

0 using the uniform distribution over $[0, T]$. For this stationary process, say $\{Q^*(t) : t \geq 0\}$,

$$\text{Cov}[Q^*(0), Q^*(t)] = \frac{1}{T} \int_0^T \text{Cov}[Q(s), Q(s+t)] ds$$

and the result follows from Lemma 9.1. It is also possible to do a direct calculation. \square

10. Conclusions

In this paper we have studied the $M_t/G/\infty$ queue with a sinusoidal arrival rate function. We obtained explicit expressions for the mean number of busy servers at time t , $m(t)$, for general service times in Theorem 4.1 and for exponential, deterministic and hyperexponential service times, respectively, in (15), (25) and (29). We also obtained explicit expressions for the peak value of m and the lag in the peak of m behind the peak of λ for general service times in (12) and (11) and for exponential and deterministic service times in (16), (18), (26) and (27).

We also investigated the performance of several approximations for $m(t)$. Assuming that we fix measuring units by letting the mean service time be one, we see that the simple stationary approximation (SSA) using the average arrival rate $\bar{\lambda}$ performs well only in the unusual case that the relative frequency γ is significantly greater than one, i.e., when there is considerably more than one cycle in the arrival process per mean service time. In this region, the pointwise stationary approximation (PSA) and the polynomial approximations do not perform well. However, PSA and the polynomial approximations improve as γ decreases, i.e., as the cycle length relative to the mean service time increases. Moreover, for γ sufficiently small (e.g., $\gamma < 1$ for exponential service times), higher order polynomial approximations provide better approximations.

Finally, we showed how general periodic arrival rate functions can be treated using Fourier series and the results for sinusoidal arrival rate functions. The case of exponential service times considered in Example 8.1 suggests that the mean function m might be well approximated by fewer terms than required for the arrival rate function λ .¹

¹ We are very grateful to Rodolfo Milito and Michael Taaffe for helpful comments.

Accepted by Linda Green; received January 24, 1991. This paper has been with the authors 2 months for 1 revision.

References

- Carrillo, M. J., "Extensions of Palm's Theorem: A Review," *Management Sci.*, 37 (1991), 739-744.
- Eick, S. G., W. A. Massey and W. Whitt, "The Physics of the $M_t/G/\infty$ Queue," *Oper. Res.* 41 (1993).
- , — and —, "Infinite-Server Approximations for Multi-Server Loss Models with Time-Dependent Arrival Rates," AT&T Bell Laboratories, Murray Hill, NJ, 1993b (in preparation).
- Green, L. and P. Kolesar, "The Pointwise Stationary Approximation for Queues with Nonstationary Arrivals," *Management Sci.*, 37 (1991), 84-97.
- , — and A. Svoronos, "Some Effects of Nonstationarity on Multiserver Markovian Queueing Systems," *Oper. Res.*, 39 (1991), 502-511.
- Heyman, D. P. and W. Whitt, "The Asymptotic Behavior of Queues with Time-Varying Arrival Rates," *J. Appl. Prob.*, 21 (1984), 143-156.
- Jagerman, D. L., "Nonstationary Blocking in Telephone Traffic," *Bell System Tech. J.*, 54 (1975), 625-661.
- Khintchine, A. Y., *Mathematical Methods in the Theory of Queueing*, Charles Griffin and Co., London, 1960 (translation of 1955 Russian book).
- Palm, C., *Intensity Variations in Telephone Traffic*, North-Holland, Amsterdam, 1988 (translation of 1943 article in *Ericsson Technics*, 44, 1-189).
- Prékopa, A., "On Secondary Processes Generated by a Random Point Distribution of Poisson Type," *Annales Univ. Sci. Budapest de Eötvös Nom. Sectio Math.*, 1 (1958), 153-170.
- Rolski, T., "Queues with Nonstationary Inputs," *Queueing Systems*, 5 (1989), 113-130.
- Ross, S. M., "Average Delay in Queues with Nonstationary Poisson Arrivals," *J. Appl. Prob.*, 15 (1978), 602-609.
- , *Stochastic Processes*, Wiley, New York, 1982.
- Rothkopf, M. H. and S. S. Oren, "A Closure Approximation for the Non-Stationary $M/M/s$ Queue," *Management Sci.*, 25 (1979), 522-534.
- Serfozo, R. F., "Point Processes," in *Handbooks in OR and MS*, Vol. 2, D. P. Heyman and M. J. Sobel, Eds., Elsevier Science Publishers, Amsterdam, 1990, 1-93.
- Thorisson, H., "On Regenerative and Ergodic Properties of the k -Server Queue with Non-Stationary Poisson Arrivals," *J. Appl. Prob.*, 22 (1985), 893-902.
- Tolstov, G. P., *Fourier Series*, Dover, New York, 1976.
- Whitt, W., "The Continuity of Queues," *Adv. Appl. Prob.*, 6 (1974), 175-183.
- , "Approximating a Point Process by a Renewal Process, I: Two Basic Methods," *Oper. Res.*, 30 (1982), 125-147.
- , "The Pointwise Stationary Approximation for $M_t/M_t/s$ Queues Is Asymptotically Correct as the Rates Increases," *Management Sci.*, 37 (1991), 307-314.
- Ziemer, R. E. and W. H. Tranter, *Principles of Communication: Systems, Modulation and Noise*, Houghton-Mifflin, Boston, MA, 1976.