# Logarithmic Asymptotics for Steady-State Tail Probabilities in a Single-Server Queue

Peter W. Glynn; Ward Whitt

# Logarithmic asymptotics for steady-state tail probabilities in a single-server queue

## PETER W. GLYNN AND WARD WHITT

### Abstract

We consider the standard single-server queue with unlimited waiting space and the first-in first-out service discipline, but without any explicit independence conditions on the interarrival and service times. We find conditions for the steady-state waiting-time distribution to have asymptotics of the form $x^{-1} \log P(W > x) \to -\theta^*$ as $x \to \infty$ for $\theta^* > 0$. We require only stationarity of the basic sequence of service times minus interarrival times and a Gärtner–Ellis condition for the cumulant generating function of the associated partial sums, i.e. $n^{-1} \log E \exp(\theta S_n) \to \psi(\theta)$ as $n \to \infty$, plus regularity conditions on the decay rate function $\psi$. The asymptotic decay rate $\theta^*$ is the root of the equation $\psi(\theta) = 0$. This result in turn implies a corresponding asymptotic result for the steady-state workload in a queue with general non-decreasing input. This asymptotic result covers the case of multiple independent sources, so that it provides additional theoretical support for a concept of effective bandwidths for admission control in multiclass queues based on asymptotic decay rates.

WAITING TIME; WORKLOAD, TAIL PROBABILITIES; ASYMPTOTICS; LARGE DEVIATIONS; CUMULANT GENERATING FUNCTION; COUNTING PROCESS; EFFECTIVE BANDWIDTHS; ADMISSION CONTROL; HIGH-SPEED NETWORKS; ASYNCHRONOUS TRANSFER MODE

AMS 1991 SUBJECT CLASSIFICATION: PRIMARY 60K25

## 1. Introduction and summary

We are pleased to be able to contribute to this Festschrift in honor of Lajos Takács on his 70th birthday. In this paper we try to emulate Takács by seeking the essential mathematics underlying a probability problem of applied relevance. Like Takács (1962), (1963), (1967), we focus on the single-server queue.

In particular, we focus on asymptotics for the steady-state waiting time $W$ and the steady-state workload $L$. We find general conditions under which

$$(1.1) \qquad\qquad x^{-1} \log P(W > x) \to -\theta^* \text{ as } x \to \infty$$

for $\theta^* > 0$, and similarly for $L$. We call the constant $\theta^*$ in (1.1) the *asymptotic decay rate*. The following elementary proposition helps put (1.1) in perspective. It is easily proved using integration by parts; see for instance p. 150 of Feller (1971).

*Proposition* 1. *For any random variable $Z$ and positive constant $\theta^*$, the following are equivalent*:

    (i) $x^{-1} \log P(Z > x) \to -\theta^*$ as $x \to \infty$;

    (ii) $\sup \{\theta \geq 0 : E \exp(\theta Z) < \infty\} = \theta^*$;

    (iii) *For all $\epsilon > 0$, there is an $x_0 \equiv x_0(\epsilon)$ such that*

$$\exp(-(\theta^* + \epsilon)x) \leq P(Z > x) \leq \exp(-(\theta^* - \epsilon)x) \text{ for all } x > x_0.$$

There is currently great interest in asymptotics such as in (1.1) because of possible applications to the design and control of emerging high-speed communication networks. In particular, it has been recognized that asymptotic decay rate functions (defined below) that determine asymptotic decay rates such as $\theta^*$ in (1.1) may be used to create effective bandwidths for admission control and other network resource allocation problems; see Gibbens and Hunt (1991), Kelly (1991), Guerin et al. (1991), Chang (1993), Chang et al. (1992), Elwalid and Mitra (1993), Sohraby (1993), Whitt (1994) and references in these sources. Chang et al. (1992) also illustrates how the asymptotic decay rates may be used to speed up simulations. Our approach here is most closely related to the papers by Whitt (1994), Chang (1993) and Chang et al. (1992). In particular, the results here provide theoretical support for the procedures in Whitt (1994).

In many cases, a stronger limit than (1.1) holds, namely,

$$(1.2) \qquad\qquad \exp(\theta^* x) P(W > x) \to \alpha^* \text{ as } x \to \infty$$

for positive constants $\theta^*$ and $\alpha^*$. Then we call $\alpha^*$ the *asymptotic constant*. It is easy to see that (1.2) implies (1.1) but not conversely. An $M/G/1$ queue for which (1.1) holds but (1.2) does *not* appears in Example 5 of Abate et al. (1994a). Then $P(W > x) \sim \alpha x^{-3/2} \exp(-\theta^* x)$ as $x \to \infty$, where $f(x) \sim g(x)$ means that $f(x)/g(x) \to 1$. *In this paper we focus on the weaker form* (1.1). For work focusing on (1.2), see Abate et al. (1994a, b, c), Asmussen (1989), (1993), Asmussen and Perry (1992), Elwalid and Mitra (1993), (1994), Neuts (1986), Stern and Elwalid (1991), Tijms (1986), van Ommeren (1988) and references in these sources.

When (1.2) holds, a natural approximation for the tail probabilities is $P(W > x) \approx \alpha^* \exp(-\theta^* x)$ for $x$ not too small. Since the asymptotic constant $\alpha^*$ in (1.2) is often not easy to obtain, Abate et al. (1994a) suggest the simple approximation $\alpha^* \approx \theta^* EW$. For some purposes, e.g. for percentiles, even $\alpha \approx 1$ is satisfactory. In many cases, $\alpha^* \approx 1$ produces a bound, i.e. $P(W > x) \leq \exp(-\theta^* x)$; see p. 269 of Asmussen (1987) and Chang (1993). These exponential approximations can also be used with (1.1), even though (1.1) does not provide as much support as (1.2). However, Example 5 of Abate et al. (1994a) shows that the quality of the approximation provided by the asymptotics can deteriorate dramatically when (1.1) holds but (1.2) does not. Moreover, for the admission control problem, it is important to note that the quality of the approximations for the tail probabilities provided by the simple one-term exponential approximations

also can deteriorate dramatically when the number of independent sources increases; see Choudhury et al. (1993), (1994).

In this first section, we present our main result and discuss its implications. We give proofs in Sections 2–8 and an example in Section 9. In Section 1.1 we state our main result for $W$; in Section 1.2 we discuss some implications and related results; in Section 1.3 we state our main results for $L$, which follow directly from the results for $W$ by discretizing the processes; and in Section 1.4 we give sufficient conditions for $W$ and $L$ to have the same logarithmic asymptotics. This involves the logarithmic asymptotics of the time-stationary and customer-stationary (embedded-stationary or Palm-stationary) versions of the arrival process. In Section 1.5 we discuss logarithmic asymptotics for steady-state queue lengths.

### 1.1. The main result

Let $\{X_n : n \geq 1\}$ be a sequence of real-valued random variables and define the associated waiting-time sequence $\{W_n : n \geq 0\}$ recursively by letting $W_0 = 0$ and

$$(1.3) \qquad W_{n+1} = [W_n + X_{n+1}]^+, \qquad n \geq 0,$$

where $[x]^+ = \max\{x, 0\}$. Let $S_0 = 0$ and $S_n = X_1 + \cdots + X_n, n \geq 1$. Let $\Rightarrow$ denote convergence in distribution.

*Theorem 1. Let $\{X_n : n \geq 1\}$ be strictly stationary. If there exists a function $\psi$ and positive constants $\theta^*$ and $\epsilon^*$ such that*

$$(1.4) \qquad \text{(i)} \quad n^{-1} \log E \exp(\theta S_n) \to \psi(\theta) \text{ as } n \to \infty \text{ for } |\theta - \theta^*| < \epsilon^*,$$

(ii) $\psi$ *is finite in a neighbourhood of* $\theta^*$

*and differentiable at* $\theta^*$ *with* $\psi(\theta^*) = 0$ *and*

$$(1.5) \qquad \psi'(\theta^*) > 0, \text{ and}$$

$$(1.6) \qquad \text{(iii)} \quad E \exp(\theta^* S_n) < \infty \text{ for } n \geq 1,$$

*then* $W_n \Rightarrow W$ *as* $n \to \infty$ *and* (1.1) *holds.*

A significant feature of Theorem 1 is that there are no independence or Markov assumptions. Instead, we have condition (1.4) involving the asymptotic behavior of the cumulant generating functions of the partial sums $S_n$, as in the Gärtner (1977)–Ellis (1984) theorem of large deviations theory; see p. 14 of Bucklew (1990). (For a discussion of the connection to cumulants, see Choudhury and Whitt (1994).) Indeed, our proof of Theorem 1 follows large deviations theory, using exponential changes of measure. For additional background on large deviations theory, see Dembo and Zeitouni (1992) and Shwartz and Weiss (1993). Theorem 3.9 (ii) of Chang (1993), which was obtained independently and concurrently, is a result closely related to Theorem 1. The upper bound was obtained earlier in Chang (1992).

The conditions in Theorem 1 are very general, but they are not necessary, as we show in Example 1 in Section 9.

A (familiar) key step in proving Theorem 1 is representing $W_n$ as the maximum of reverse-time partial sums; i.e.

$$(1.7) \qquad W_n = S_n - \min_{0 \le k \le n} S_k = \max_{0 \le k \le n} \{S_n - S_k\}.$$

so that, when we extend $\{X_n\}$ to a doubly infinite stationary sequence $\{X_n : -\infty < n < \infty\}$, $W_n$ is distributed as $\tilde{M}_n \equiv \max\{\tilde{S}_k : 0 \le k \le n\}$, where $\tilde{S}_0 = 0$ and $\tilde{S}_k = X_{-1} + \cdots + X_{-k}$. The conditions in Theorem 1 obviously apply to $\tilde{S}_k$ as well as $S_k$ because $E \exp(\theta \tilde{S}_k) = E \exp(\theta S_k)$. Since the stationarity is required only for this step, we obtain Theorem 1 immediately from the following result for maxima of partial sums $M_n = \max\{S_k : 0 \le k \le n\}$, which does not require stationarity. We prove the following result in Section 2.

*Theorem 2. Let $\{X_n : n \ge 1\}$ be a sequence of real-valued random variables (not necessarily stationary or mutually independent). If there exists a function $\psi$ and positive constants $\theta^*$ and $\epsilon^*$ such that (1.4)–(1.6) hold and*

$$(1.8) \qquad \limsup_{n \to \infty} E \exp(\theta X_n) < \infty \text{ for } |\theta| < \epsilon^*,$$

*then $M_n \to M$ w.p.1 as $n \to \infty$ and*

$$(1.9) \qquad x^{-1} \log P(M > x) \to -\theta^* \text{ as } x \to \infty.$$

Note that condition (1.8) in Theorem 2 is not needed if $\{X_n\}$ is stationary, because then (1.8) is implied by (1.6) in the case $n = 1$. For this, recall that $E \exp(\theta_1 Z) < \infty$ when $E \exp(\theta_2 Z) < \infty$ and $\theta_1 < \theta_2$ for any random variable $Z$ by Hölder's inequality; see (21) on p. 47 of Chung (1974).

Also note that condition (1.6) is clearly necessary in Theorem 2, because $M \ge S_n$ for all $n$. Hence, $E \exp(\theta^* M) = \infty$ if $E \exp(\theta^* S_n) = \infty$ for any $n$.

We remark that we have also proved a version of Theorem 2 with condition (1.8) replaced by $E \exp(\theta S_n) < \infty$ for $n \ge 1$ for some $\theta$ with $\theta > \theta^*$. This alternative condition might be preferred in Theorem 2, but it would require that we strengthen (1.6) in Theorem 1.

In Theorem 1 we have assumed that the basic sequence $\{X_n\}$ is stationary. However, this is not a great restriction because the focus is on the steady-state waiting time $W$. Given the distribution of $W$, it is usually possible to choose a stationary version of any given basic sequence $\{X_n\}$ such that $W_n \Rightarrow W$ as $n \to \infty$; see for example p. 13 of Borovkov (1976). Of course, the conditions in Theorem 1 apply to this stationary version. However, under regularity conditions, non-stationary versions and stationary versions of the basic sequence will couple so that the conditions for one enable us to verify the conditions for the other.

In other words, $W$ typically does not depend on the initial part of the basic sequence $\{X_n\}$. In contrast, the maximum $M$ in Theorem 2 clearly does depend on

the entire sequence $\{X_n\}$. For a simple example, suppose that $\{X_n : n \geq 2\}$ is i.i.d. with a good distribution, but $P(X_1 > x) \sim Ax^{-p}$. Then $X_1$ influences the tail behavior of $S_n$ for all $n$ and $x$, but not that of $W$.

An important role is played by the function $\psi$ in Theorem 1; we call it the *(asymptotic) decay rate function*. It is significant that $\psi$ is necessarily convex where it is finite, because $\log E \exp(\theta Z)$ is convex where it is finite for any random variable $Z$, as can easily be seen by applying Hölder's inequality. It is important to distinguish the decay rate function $\psi$ from the associated *large deviations rate function $I(x)$*, defined by

$$(1.10) \qquad\qquad I(x) = \sup_\theta \{\theta x - \psi(\theta)\};$$

see for example Chapter 1 of Bucklew (1990). The functions $\phi$ and $I$ are intimately related. Indeed, they are convex conjugates of each other; see p. 183 of Bucklew (1990).

### 1.2. *Implications and related results*

The conditions of Theorems 1 and 2 are easy to check when the basic sequence $\{X_n\}$ is i.i.d. This special case includes the $GI/GI/1$ queue (with i.i.d. service times independent of i.i.d. interarrival times), for which it is possible to obtain the stronger result (1.2); see for example p. 269 of Asmussen (1987). For early results in this direction, see Smith (1953) and Theorems V.10.1 and VI.6.1 of Keilson (1965). Asmussen (1987), (1993) refers to the history in risk theory. In this $GI/GI/1$ case, Abate et al. (1993) have shown that it is also easy to compute the tail probabilities by numerical transform inversion, numerically integrating a contour integral representation for $E \exp(-\theta W)$.

*Corollary* 1. *If $\{X_n : n \geq 1\}$ is i.i.d., $EX < 0$,*

$$(1.11) \qquad\qquad E \exp(\theta^* X) = 1$$

*and $E \exp(\theta X) < \infty$ for $-\epsilon < \theta < \theta^* + \epsilon$ for some $\epsilon > 0$, then the conditions of Theorems 1 and 2 hold with $\psi(\theta) = \log E \exp(\theta X)$, so that (1.1) and (1.9) hold.*

*Proof.* Note that $n^{-1} \log E \exp(\theta S_n) = E \exp(\theta X_1) = \psi(\theta)$ when $\{X_n\}$ is i.i.d. Since $\psi(0) = 1$, $\psi'(0) = EX < 0$ and $\psi$ is convex, $\psi'(\theta^*) > 0$.

Thus, for the $GI/GI/1$ queue, it is easy to see what the decay rate function $\psi$ is. For example, in the $M/M/1$ queue with service rate 1 and arrival rate $\rho$, $\psi(\theta) = -\log[(1 - \theta)(1 + \theta/\rho)]$; for the $D/M/1$ queue, $\psi(\theta) = -\log(1 - \theta) - \theta/\rho$; and for the $M/D/1$ queue, $\psi(\theta) = \theta - \log(1 + \theta/\rho)$.

It is worth pointing out that the logarithmic asymptotics in (1.1) tend to be robust. In general, weak convergence of distributions does not imply that large deviations asymptotics converge. However, in this context, weak convergence plus uniform integrability does imply that the cumulant generating function converges, and the logarithmic asymptotics here depend only on the location of the root (and

not, for example, the slope at the root). We illustrate by stating a concrete result in the context of Corollary 1.

*Corollary* 2. *Let* $\{X_n^\gamma : n \geq 1\}$ *be i.i.d. for each* $\gamma > 0$ *where* $X_1^\gamma \Rightarrow X_1^0$ *as* $\gamma \to 0$ *and* $E \exp(\theta X_1^\gamma) < M$ *for some* $\theta > \theta_0^*$ *and some finite* $M$, *for all* $\gamma$. *If* $X_1^\gamma$ *satisfies the conditions of Corollary* 1 *for each* $\gamma \geq 0$, *then* (1.1) *holds for each* $\gamma \geq 0$ *and* $\theta_\gamma^* \to \theta_0^*$ *as* $\gamma \to 0$.

*Proof.* Since $X_1^\gamma \Rightarrow X_1^0$ as $\gamma \to 0$, $\exp(\theta X_1^\gamma) \Rightarrow \exp(\theta X_1^0)$ as $\gamma \to 0$. The uniform moment bound implies the uniform integrability needed to obtain $E \exp(\theta X_1^\gamma) \to E \exp(\theta X_1^0)$ as $\gamma \to 0$ for all $\theta < \theta_0^* + \epsilon$ for some $\epsilon$.

In order to understand what the asymptotic decay rate $\theta^*$ in (1.1) primarily depends upon, and sometimes to compute $\theta^*$, it is useful to consider heavy-traffic asymptotic expansions for $\theta^*$ based on expanding the function $\psi(\theta)$ in a Taylor series expansion about 0. Such heavy-traffic asymptotic expansions are established in Abate et al. (1994a), Abate and Whitt (1994) and Choudhury and Whitt (1994). Since $\log E \exp(\theta S_n)$ is the cumulant generating function of $S_n$, the derivatives of $\psi(\theta)$ are the asymptotic cumulants of $S_n$. To illustrate, we establish the first term of the heavy-traffic expansion here. The first term coincides with the familiar decay rate associated with exponential heavy-traffic limits; see Kingman (1962), Iglehart and Whitt (1970) and Choudhury and Whitt (1994).

*Corollary* 3. *Consider a family of models indexed by* $\rho$, $0 < \rho < 1$. *Suppose that the assumptions of Theorem* 1 *hold for each* $\rho$ *and*
(i)   $EX_n(\rho) = -(1 - \rho)$,
(ii)  $n^{-1} \operatorname{Var} S_n(\rho) \to \sigma^2$ *as* $n \to \infty$, $0 < \sigma^2 < \infty$, *and*
(iii) $n^{-1} E(S_n(\rho) - (1 - \rho)n)^3 \to \gamma$ *as* $n \to \infty$, $-\infty < \gamma < \infty$.
*Then* (1.1) *holds with*

$$\theta^*(\rho) = \frac{2(1 - \rho)}{\sigma^2} + o(1 - \rho) \ as \ \rho \to 1.$$

*Proof.* Since $\log E \exp(\theta S_n)$ is the cumulant generating function of $S_n$, we can apply Taylor's theorem to obtain

$$n^{-1} \log E \exp(\theta S_n) = -\theta(1 - \rho) + \frac{\theta^2 \operatorname{Var}(S_n)}{2n} + o(\theta^2) \ as \ \theta \to 0$$

uniformly in $\rho$ and $\eta$, using condition (iii) to get the uniformity in $n$; e.g., see $(4')$ on p. 268 of Chung (1974). Hence

$$\psi(\theta) = -\theta(1 - \rho) + \frac{\theta^2 \sigma^2}{2} + o(\theta^2) \ as \ \theta \to 0$$

uniformly in $\rho$, so that the desired conclusion follows.

Another easy case is when the partial sums $S_n$ are Gaussian (but possibly dependent) for all $n$. When $S_n$ is Gaussian, Theorem 1 takes a very simple form. In particular, then (1.1) holds with $\theta^*$ in Corollary 3. The Gaussian assumption holds approximately in an $E_k/E_m/1$ queue for suitably large $k$ and $m$. (As usual, $E_k$ stands for Erlang of order $k$.) A direct Gaussian approximation has also been proposed and studied by Addie and Zuckerman (1993). This analysis provides additional justification for the heavy-traffic approximation, because it does not (at least directly) require a high traffic intensity.

*Corollary 4. Suppose that $S_n$ is Gaussian with negative mean $m_n$ and finite variance $\sigma_n^2$ for all $n \geq 1$. If $m_n/n \to m$ and $\sigma_n^2/n \to \sigma^2$ as $n \to \infty$, where $m < 0 < \sigma^2$, then the conditions of Theorem 2 hold with $\psi(\theta) = \theta m + \theta^2 \sigma^2/2$, so that (1.9) holds with $\theta^* = -2m/\sigma^2 > 0$. If, in addition, the basic sequence $\{X_n\}$ is stationary, then the conditions of Theorem 1 hold, so that (1.1) holds.*

*Proof.* Recall that $E \exp(\theta S_n) = \exp(\theta m_n + \theta^2 \sigma^2/2)$ when $S_n$ is Gaussian with mean $m_n$ and variance $\sigma_n^2$.

In queueing theory, (1.3) is the familiar Lindley equation associated with a single-server queue with unlimited waiting room and the first-in first-out service discipline. Then $X_n = V_n - U_n$ where, for $n \geq 1$, $V_n$ is the service time of customer $n$ and $U_n$ is the interarrival time between customers $n$ and $n + 1$. With this indexing, we begin with a first customer arriving at an empty system.

Another queueing model that leads to the representation $X_n = V_n - U_n$ is the queue length in a discrete-time single-server queue. Then we interpret $V_n$ as the number of arrivals at epoch $n$ and $U_n$ as the number of *potential* departures at epoch $n$. For this representation to be valid, we usually require special Markov or deterministic assumptions in the service process, or 'autonomous service'; see p. 235 of Borovkov (1976). We use this below in Section 1.3. For the ATM networks it is often reasonable to assume deterministic service, so that this Lindley equation representation is indeed appropriate. For example, if there is at most one service completion at each epoch, then $U_n = 1$ for all $n$. This model variant is considered by Chang (1993), Sohraby (1993) and Chang et al. (1992). Chang (1993) also focuses on the Gärtner–Ellis condition in (1.4).

Given Theorems 1 and 2, we want to know when the conditions are satisfied. In the queueing context, the conditions can be expressed in terms of the two sequences $\{U_n : n \geq 1\}$ and $\{V_n : n \geq 1\}$ separately when the sequences $\{U_n\}$ and $\{V_n\}$ are independent. (However, note that such independence is not required in Theorem 1.)

To state the result, let $S_n^v = V_1 + \cdots + V_n$ and let $S_n^u = U_1 + \cdots + U_n$.

*Proposition 2. Suppose that $X_n = V_n - U_n$, where $\{V_n : n \geq 1\}$ and $\{U_n : n \geq 1\}$ are independent sequences of non-negative random variables (without independence or stationarity assumptions for each sequence). If there exist functions $\psi_v$ and $\psi_u$*

*and positive constants $\theta^*$, $\epsilon^*$ and $M$ such that*

(1.12)   (i)     $n^{-1} \log E \exp(\theta S_n^v) \to \psi_v(\theta)$ *as $n \to \infty$ for $|\theta - \theta^*| < \epsilon^*$,*

(1.13)   (ii)    $\psi_v$ *is finite in a neighborhood of $\theta^*$ and differentiable at $\theta^*$,*

(1.14)   (iii)   $E \exp(\theta^* S_n^v) < \infty$ *for $n \geq 1$,*

(1.15)   (iv)    $E \exp(\theta V_n) < M$ *for $n \geq 1$ and all $|\theta| < \epsilon^*$,*

(1.16)   (v)     $n^{-1} \log E \exp(-\theta S_n^u) \to \psi_u(-\theta)$ *as $n \to \infty$, for $|\theta - \theta^*| < \epsilon^*$,*

(1.17)   (vi)    $\psi_u$ *is finite in a neighborhood of $-\theta^*$ and differentiable at $-\theta^*$,*

(1.18)   (vii)   $E \exp(-\theta U_n) < M$ *for $n \geq 1$ and all $|\theta| < \epsilon^*$ and*

(1.19)   (viii)  $\psi(\theta^*) = 0$ *and $\psi'(\theta^*) > 0$ for $\psi(\theta) = \psi_v(\theta) + \psi_u(-\theta)$,*

*then $\{X_n\}$ satisfies conditions (1.4)–(1.6) and (1.8) with decay rate function $\psi$, so that (1.9) holds. If, in addition, $\{(U_n, V_n)\}$ is stationary, then $\{X_n\}$ is stationary and (1.1) holds.*

*Proof.* By the independence,

$$E \exp(\theta S_n) = E \exp(\theta(S_n^v - S_n^u)) = E \exp(\theta S_n^v) E \exp(-\theta S_n^u),$$

so that

$$\log E \exp(\theta S_n) = \log E \exp(\theta S_n^v) + \log E \exp(-\theta S_n^u).$$

Since $S_n^u \geq 0, E \exp(-\theta S_n^u) \leq 1$. Similarly, $E \exp(\theta X_n) = E \exp(\theta V_n) E \exp(-\theta U_n)$. Hence, it is clear that the assumed conditions here imply the conditions in Theorem 1.

Assuming that the arrival and service processes are independent, Proposition 2 shows that we can treat them separately, in the sense that the overall decay rate function $\psi$ is the sum of the component decay rate functions, as indicated in (1.19). This separability is a basic feature of the heavy-traffic limits in Iglehart and Whitt (1970) and in the effective bandwidth approximations; see Elwalid and Mitra (1993), (1994), Stern and Elwalid (1991) and Whitt (1994). To obtain further results for these separate processes, it is useful to have a relation between the asymptotics for a counting process and the asymptotics for its inverse partial sum process. For this purpose, we apply a result from Glynn and Whitt (1994).

Let $\{T_n : n \geq 0\}$ be a non-decreasing sequence of random variables with $T_0 = 0$. We think of $T_n$ as the arrival epoch of customer $n$ in the queue; then $T_n = U_1 + \cdots + U_n$. Let $\{N(t) : t \geq 0\}$ be the associated counting process defined by

(1.20)                     $N(t) = \max\{n \geq 0 : T_n \leq t\}, \qquad t \geq 0.$

The (familiar) key relation between $T_n$ and $N(t)$ that we exploit is

$$(1.21) \qquad \{N(t) \geq n\} = \{T_n \leq t\}$$

for all non-negative $n$ and $t$.

A process $\{Z(t) : t \geq 0\}$ will be said to satisfy the *Gärtner–Ellis condition with decay rate function* $\psi$ if

$$(1.22) \qquad \lim_{t \to \infty} t^{-1} \log E \exp^{\theta Z(t)} = \psi(\theta) \text{ for all } \theta \in \mathbb{R}.$$

For a discrete-time process, we let $t$ run through the positive integers in (1.22).

The associated decay rate function $\psi$ will be said to satisfy the *auxiliary large deviations (LD) regularity conditions* if (1.23)–(1.26) below hold:

$$(1.23) \qquad \beta \equiv \inf \{\theta : \psi(\theta) = +\infty\} > 0,$$

$$(1.24) \qquad \psi \text{ is differentiable everywhere in } (-\infty, \beta),$$

$$(1.25) \qquad \lim_{\theta \uparrow \beta} \psi'(\theta) = +\infty,$$

and

$$(1.26) \qquad \lim_{\theta \uparrow \beta} \psi(\theta) = \psi(\beta).$$

The conditions (1.22)–(1.26) are standard in the large deviations literature. In particular, under conditions (1.22)–(1.26), the process $\{Z(t) : t \geq 0\}$ satisfies the Gärtner (1977)–Ellis (1984) theorem, i.e. the *large deviations principle* holds for $\{Z(t) : t \geq 0\}$ with rate function I in (1.10); see pp. 42–50 of Dembo and Zeitouni (1992).

The following result is proved in Glynn and Whitt (1994). Let $\psi^{-1}$ be the inverse function of $\psi$. Note that $\psi$ is non-decreasing, and strictly increasing where it is finite. Hence, for $x$ and $y$ finite, $\psi^{-1}(y) = x$ if and only if $\psi(x) = y$.

*Theorem 3. If the counting process $\{N(t) : t \geq 0\}$ satisfies (1.22)–(1.26), then the inverse partial sum process $\{T_n : n \geq 0\}$ does too, with the possible exception of (1.22) for $\theta = \beta_T$. Similarly, if $\{T_n : n \geq 0\}$ satisfies (1.22)–(1.26), then $\{N(t) : t \geq 0\}$ does too, with the possible exception of (1.22) for $\theta = \beta_N$. In particular, then (1.22) holds for both processes, i.e.*

$$(1.27) \qquad \lim_{t \to \infty} t^{-1} \log E \exp(\theta N(t)) = \psi_N(\theta)$$

*and*

$$(1.28) \qquad \lim_{n \to \infty} n^{-1} \log E \exp(\theta T_n) = \psi_T(\theta)$$

*both hold (with the noted exceptions) and*

$$(1.29) \qquad \psi_N(\theta) = -\psi_T^{-1}(-\theta),$$

*where they are finite.*

Thus, subject to regularity conditions, given the Gärtner–Ellis asymptotics for one of $N$ or $T$, we obtain the Gärtner–Ellis asymptotics for the other directly and have the inverse relation (1.29). This parallels previous relations between other limits for $N$ and $T$; see for example Iglehart and Whitt (1971), Section 7 of Whitt (1980), Theorems 3 and 6 of Glynn and Whitt (1988a) and Theorem 1 of Glynn and Whitt (1988b).

For example, we can apply Theorem 3 to obtain the Gärtner–Ellis limit (1.28) for the partial sums $S_n^u$ from the Gärtner–Ellis limit (1.27) for the counting process $N(t)$ derived for batch Markovian arrival processes in Theorem 1 of Choudhury and Whitt (1994). Abate et al. (1994c) obtain (1.2) for $BMAP/GI/1$ queues, while the results here yield (1.1) for $BMAP/G/1$ queues, without requiring that the service times be i.i.d. Sufficient conditions for (1.22) in terms of embedded regenerative structure are also given in Theorem 7 of Glynn and Whitt (1994).

We now show that deterministic sequences provide upper bounds on $\theta^*$ when $\{U_n\}$ and $\{V_n\}$ are independent sequences; see Section 8 of Abate et al. (1994a) for related results. We use the queueing notation $G/G/1$ to refer to a general stationary sequence $\{(U_n, V_n)\}$ of interarrival times and service times.

*Proposition 3. Among $G/G/1$ models satisfying the assumptions of Proposition 2, the asymptotic decay rate $\theta^*$ is maximized (a) by deterministic service times among all stationary service-time sequences $\{V_n\}$ with given mean $EV_n$, and (b) by deterministic interarrival times among all stationary interarrival-times sequences $\{U_n\}$ with given mean $EU_n$.*

*Proof.* By Jensen's inequality, $E \exp(\theta Z) \geq \exp(\theta E Z)$ for any random variable $Z$, so that $\log E \exp(\theta S_n^v) \geq \log \exp(\theta E S_n^v) = n\theta E V_1$ and $\log E \exp(-\theta S_n^u) \geq \log \exp(\theta E S_n^u) = n\theta E U_1$. Hence, if $\psi_v^D$ and $\psi_u^D$ denote the decay rate functions in the deterministic cases, then $\psi_v(\theta) \geq \psi_v^D(\theta)$ and $\psi_u(-\theta) \geq \psi_u^D(-\theta)$ for all $\theta > 0$, so that the roots in (1.19) must be ordered as indicated.

More generally, we can establish stochastic comparisons between any two $G/G/1$ systems.

*Proposition 4. Consider two $G/G/1$ queues satisfying the assumptions of Theorem 1. If $E \exp(\theta S_n^1) \leq E \exp(\theta S_n^2)$ for all $\theta \geq 0$ and all $n$ suitably large, then $\theta_2^* \leq \theta_1^*$.*

*Proof.* The condition implies that $\psi_1(\theta) \leq \psi_2(\theta)$ for all $\theta \geq 0$. Hence, the roots $\theta_i^*$ of $\psi_i(\theta) = 0$ must be ordered by $\theta_2^* \leq \theta_1^*$.

As in Whitt (1994), when $\{U_n\}$ and $\{V_n\}$ are independent, we can characterize the arrival and service decay rate functions $\psi_u(-\theta)$ and $\psi_v(\theta)$ from the asymptotic decay rates $\theta^*$ observed in $G/D/1$ and $D/G/1$ queues. To do this, we must consider all possible arrival rates $\rho$, $0 < \rho < 1$, so that the asymptotic decay rate $\theta^*$ becomes a *function* $\theta^*(\rho)$, $0 < \rho < 1$. Let $\psi_u(-\theta)$ refer to the case in which

$EU_n = 1$ and let the case of arrival rate $\rho$ be obtained by considering interarrival times $U_n/\rho$ for all $n$, i.e. simple time scaling.

**Proposition 5.** *For $G/G/1$ models satisfying the assumptions of Proposition 2,*

(i) *the arrival asymptotic decay rate function $\psi_u(-\theta)$ with arrival rate 1 is determined by the decay rate $\theta^*(\rho)$ in $G/D/1$ models with arrival rate $\rho$, $0 < \rho < 1$, i.e. by the equation*

$$(1.30) \qquad \psi_u(-\theta^*(\rho)/\rho) + \theta^*(\rho) = 0, \qquad 0 < \rho < 1.$$

(ii) *The service asymptotic decay rate function $\psi_v(\theta)$ with service rate 1 is determined by the decay rate $\theta^*(\rho)$ in $D/G/1$ models with arrival rate $\rho$, $0 < \rho < 1$, i.e. by the equation*

$$(1.31) \qquad \psi_v(\theta^*(\rho)) - \theta^*(\rho)/\rho = 0.$$

*Proof.* Note that $\psi_u(-\theta)$ is a decreasing convex function with $\psi_u'(0) = -1$. Hence, the values of $\psi_u(-\theta)$ for $\theta > 0$ are determined by the intersection with all lines through the origin with slopes less than $-1$. This is provided by (1.30), after making the change of variables $\theta(\rho) = \theta^*(\rho)/\rho$. Similarly, $\psi_v(\theta)$ is an increasing convex function with $\psi_v'(0) = 1$. Hence, the values of $\psi_v(\theta)$ for $\theta > 0$ are determined by the intersection with all lines through the origin with slope greater than $+1$. This is determined by (1.31).

Since Chang's (1993) model is equivalent to the $D/G/1$ special case, his equation $a^*(\theta) = c$ in (64) and the proof of Theorem 3.9 (ii) should coincide with (1.31), and it does.

### 1.3. *A continuous-time analog: the workload*

We can apply Theorems 1 and 2 to obtain corresponding results for continuous-time workload processes; we will only discuss the analog of Theorem 1. Paralleling (1.7), suppose that we have a continuous-time workload process $\{L(t) : t \geq 0\}$ defined in terms of a continuous-time net input process $\{Y(t) : t \geq 0\}$ by applying the usual reflection map, i.e.

$$(1.32) \qquad L(t) = Y(t) - \inf\{Y(s) : 0 \leq s \leq t\}, \qquad t \geq 0,$$

with $L(0) = 0$. Moreover, let the net input process be defined in terms of a total input process $\{I(t) : t \geq 0\}$ with non-decreasing sample paths by

$$(1.33) \qquad Y(t) = I(t) - t, \qquad t \geq 0.$$

In the $G/G/1$ queue, $I(t)$ represents the total work in service time to arrive in the interval $[0, t]$, i.e. the sum of all service times of all arrivals in $[0, t]$, but here $I(t)$ can be more general. For example, this formulation includes fluid models such as the Markov modulated fluid models in Elwalid and Mitra (1993) as a special case (without directly requiring the Markov assumption).

Paralleling Theorem 1, we will work with a version of $I(t)$ that has stationary increments. We prove the following result in Section 4.

*Theorem 4. Let the net input process $\{Y(t) : t \geq 0\}$ have stationary increments with $EY(t) = (\rho - 1)t$ where $\rho < 1$. If there exist a function $\psi$ and positive constants $\theta^*$ and $\epsilon^*$ such that the analogs of (1.4) and (1.6) hold, i.e. if*

$$(1.34) \qquad t^{-1} \log E \exp\left(\theta Y(t)\right) \to \psi(\theta) \text{ as } t \to \infty \text{ for } |\theta - \theta^*| < \epsilon^*$$

*and*

$$(1.35) \qquad\qquad\qquad E \exp\left(\theta^* Y(t)\right) < \infty \text{ for all } t > 0,$$

*and if (1.5) holds for this $\psi$, then $L(t) \Rightarrow L$ as $t \to \infty$ and*

$$(1.36) \qquad\qquad x^{-1} \log P(L > x) \to -\theta^* \text{ as } x \to \infty.$$

Theorem 4 easily applies to superpositions of independent processes, as we now show.

*Proposition 6. Consider the workload process $L(t)$ in (1.32) and (1.33) with $I(t) = I_1(t) + \cdots + I_n(t)$, where $I_1(t), \ldots, I_n(t)$ are mutually independent nondecreasing processes each with stationary increments satisfying*

$$t^{-1} \log E \exp\left(\theta I_i(t)\right) \to \psi_i(\theta) \text{ for } |\theta - \theta^*| < \epsilon^*$$

*and*

$$E \exp\left(\theta^* I_i(t)\right) < \infty \text{ for all } t > 0 \text{ and } i.$$

*If (1.5) holds for*

$$\psi(\theta) = \psi_1(\theta) + \cdots + \psi_n(\theta) - \theta,$$

*then the conditions of Theorem 4 hold, so that (1.36) holds.*

*Proof.* By the independence,

$$\log E \exp\left(\theta Y(t)\right) = \log E \exp\left(\theta I_1(t)\right) + \cdots + \log E \exp\left(\theta I_n(t)\right) - \theta t.$$

The following proposition treats the standard case in queueing, in which the total input $I(t)$ is the sum of all the service times of all arrivals in the interval $[0, t]$. We prove the following in Section 7.

*Proposition 7. Consider a total input process defined by*

$$(1.37) \qquad\qquad I(t) = \sum_{i=1}^{A(t)} V_i, \qquad t \geq 0.$$

*Suppose that $\{V_n\}$ is independent of $\{A(t)\}$,*

$$n^{-1} \log E \exp\left(\theta \sum_{i=1}^{n} V_i\right) \to \psi_v(\theta) \text{ as } n \to \infty \text{ for all } \theta \text{ in a neighborhood of } \hat{\theta}$$

*and*

$$t^{-1} \log E \exp (\theta A(t)) \to \psi_A(\theta) \text{ as } t \to \infty \text{ for all } \theta \text{ in a neighborhood of } \psi(\hat{\theta}),$$

*where $\psi_A$ is continuous at $\psi_v(\hat{\theta})$. Then*

$$t^{-1} \log E \exp (\hat{\theta} I(t)) \to \psi_A(\psi_v(\hat{\theta})) \text{ as } t \to \infty.$$

Results related to Propositions 6 and 7 also appear in Chang (1993).

### 1.4. *Palm equivalence for the Gärtner–Ellis limits: relating W and L*

The asymptotics for $W$ and $L$ differ, in part, because $W$ is based on the customer-stationary (embedded-stationary or Palm-stationary) sequence $\{U_n\}$ while $L$ is based on the counting process $\{A(t)\}$ with stationary increments, which is associated with the time-stationary sequence, say $\{U_n^*\}$, connected by the Palm transformation, see for example Franken et al. (1981). However, we anticipate that we should have $\theta_W^* = \theta_L^*$. To establish that relation, we would like to have *Palm equivalence for the Gärtner–Ellis limits*, i.e. we would like to be able to say that $n^{-1} \log E \exp (\theta S_n^u) \to \psi_u(\theta)$ as $n \to \infty$ if and only if $n^{-1} \log E \exp (\theta S_n^{u*}) \to \psi_u^*(\theta)$ as $n \to \infty$ and $\psi_u = \psi_u^*$, where $S_n^{u*} = U_1^* + \cdots + U_n^*$, $n \geq 1$. We establish a weaker result here. We show that if both limits hold with the limit functions $\psi_u$ and $\psi_u^*$ satisfying regularity conditions, then $\psi_u = \psi_u^*$. We then use this property to provide conditions under which $\theta_L^* = \theta_W^*$.

We start by relating the asymptotics for $L$ and $W$ when the service times are i.i.d. As in Theorem 2 and Section 2 of Abate et al. (1994b), we apply the generalized Takács (1963) relation between $W$ and $L$ in a $G/GI/1$ queue (with i.i.d. service times that are independent of the arrival process); see (1.38) below and (4.5.9) on p. 129 of Franken et al. (1981).

**Proposition 8.** *In a $G/GI/1$ queue, (1.1) holds if and only if (1.36) holds and $\theta_W^* = \theta_L^*$.*

*Proof.* The generalized Takács relation yields

$$(1.38) \qquad E \exp (\theta L) = 1 - \rho + \rho E \exp (\theta W) E \exp (\theta V_e),$$

where $V_e$ has the stationary-excess or equilibrium-residual-life distribution associated with the service-time distribution. By Proposition 3 here, Theorem 10 of Abate et al. (1994a) and Lemma 1 of Abate et al. (1994b), $E \exp (\theta V_e) < \infty$ if $E \exp (\theta W) < \infty$.

We now apply Proposition 8 to obtain a form of Palm equivalence for the Gärtner–Ellis limits. We prove the following result in Section 5.

**Theorem 5.** *Let $A(t)$ be a counting process associated with a non-deterministic time-stationary sequence $\{U_n^*\}$ and let $S_n^u$ be the partial sums associated with the corresponding customer-stationary sequence $\{U_n\}$. Assume that the decay rate*

*function associated with $A$, $\psi_A$, satisfies the auxiliary LD regularity conditions (1.23)–(1.26) with limit of support $\beta_A$ in (1.23). Assume that*

$$(1.39) \qquad\qquad E \exp\left(\theta A(t)\right) < \infty \text{ for all } t > 0 \text{ and } \theta < \beta_A$$

*and*

$$(1.40) \qquad\qquad t^{-1} \log E \exp\left(\theta A(t)\right) \to \psi_A(\theta) \text{ as } t \to \infty \text{ for } \theta < \beta_A.$$

*Assume that*

$$(1.41) \qquad n^{-1} \log E \exp\left(-\theta S_n^u\right) \to \psi_u(-\theta) \text{ as } n \to \infty \text{ for all } \theta > 0,$$

*where $\psi_u(-\theta)$ is finite and differentiable for all $\theta > 0$. Then*

$$(1.42) \qquad\qquad \psi_u(-\theta) = \psi_u^*(-\theta) = -\psi_A^{-1}(\theta) \text{ for all } \theta > 0.$$

We now relate the logarithmic asymptotics for $W$ and $L$ in a general $G/G/1$ queue when the arrival and service processes are independent (but the service times need not be i.i.d.). We prove the following result in Section 6.

*Theorem 6. Consider a $G/G/1$ queue in which the service times $\{V_n\}$ are stationary and independent of the arrival process. Let the arrival process satisfy the assumptions of Theorem 5. Let the service decay rate function $\psi_v$ satisfy the auxiliary LD regularity conditions (1.23)–(1.26) with limit of support $\beta_v$ in (1.23). Assume that*

$$(1.43) \qquad\qquad E \exp\left(\theta S_n^v\right) < \infty \text{ for all } n \geq 1 \text{ and } \theta < \beta_v$$

*and*

$$(1.44) \qquad n^{-1} \log E \exp\left(\theta S_n^v\right) \to \psi_v(\theta) \text{ as } n \to \infty \text{ for } \theta < \beta_v.$$

*Then (1.1) and (1.36) both hold with $\theta_W^*(\rho) = \theta_L^*(\rho)$ for each $\rho$, $0 < \rho < 1$.*

### 1.5. Queue lengths

In this section we discuss the logarithmic asymptotics for the steady-state queue length (number in system). Let $Q$ and $Q^a$ be the steady-state queue length at an arbitrary time and at an arrival epoch, respectively, which we assume are well defined. As in Section 1.3, let $G/GI/1$ mean i.i.d. service times that are independent of general stationary interarrival times. We prove the following in Section 8.

*Proposition 9. In the $G/GI/1$ queue, (1.1) holds if and only if the analogs of (1.1) hold for $Q$ and $Q^a$, in which case*

$$(1.45) \qquad\qquad \theta_Q^* = \theta_{Q^a}^* = \log E \exp\left(\theta_W^* V_1\right) = \psi_v(\theta_W^*).$$

## 2. Proof of Theorem 2

In this section we prove Theorem 2. For this purpose, we perform a change of measure for each $n$. In particular, for each $n \geq 1$, let $P_n^*$ be the probability measure on $\mathbb{R}^n$ defined by

$$P_n^*(dx_1, \ldots, dx_n) = \frac{\exp{(\theta^* S_n)}}{E \exp{(\theta^* S_n)}} P(X_1 \in dx_1, \ldots, X_n \in dx_n)$$

(2.1)
$$= \exp{(\theta^* S_n - \psi_n(\theta^*))} P(X_1 \in dx_1, \ldots, X_n \in dx_n),$$

where $\psi_n(\theta) = \log E \exp{(\theta S_n)}$ and $\psi_n(\theta^*) < \infty$ for $n \geq 1$ by (1.6).

We base our proof on the following strengthened form of the weak law of large numbers. This is closely related to claim 1 on p. 17 of Bucklew (1990) in his proof of the Gärtner–Ellis theorem. However, we only make assumptions locally around $\theta^*$, whereas Bucklew's assumptions are more global. We will need the cases $k = 0$ and $k = 1$ in our proof of Theorem 2. We prove Theorem 7 in Section 3.

*Theorem 7. Let $k$ be a fixed non-negative integer. Under the conditions of Theorem 2 (excluding (1.8) if $k = 0$), for each $\epsilon > 0$ there exists $n_0$ and $\eta \equiv \eta(\epsilon) \in [0, 1)$ such that*

(2.2)
$$P_n^* \left( \left| \frac{S_{n-k}}{n} - \nu \right| > \epsilon \right) \leq \eta^n \text{ for } n \geq n_0.$$

Since $M_n$ is non-decreasing, $M_n \to M$ w.p.1. The desired result (1.9) implies that $M$ must be proper. Since $P(M > x) = P(T(x) < \infty)$, where

(2.3)
$$T(x) = \inf{\{n \geq 0 : S_n > x\}},$$

it suffices to show that

(2.4)
$$x^{-1} \log P(T(x) < \infty) \to -\theta^* \text{ as } x \to \infty.$$

Let $\lfloor x \rfloor$ be the greatest integer less than or equal to $x$ and let $\lceil x \rceil$ be the least integer greater than or equal to $x$. Now, for $\epsilon$ and $\nu$ given, and any $x$ and $n(\epsilon)$,

$$P(T(x) < \infty) = \sum_{j=1}^{\infty} P(T(x) = j)$$

$$\leq \sum_{j=1}^{n(\epsilon)} P(T(x) = j) + \sum_{j=n(\epsilon)+1}^{\lfloor x(1-\epsilon)/\nu \rfloor} P(T(x) = j)$$

(2.5)
$$+ \sum_{j=\lceil x(1-\epsilon)/\nu \rceil}^{\lfloor x(1+\epsilon)/\nu \rfloor} P(T(x) = j) + \sum_{j=\lceil x(1+\epsilon)/\nu \rceil}^{\infty} P(T(x) = j).$$

Given $\epsilon$, we choose $n(\epsilon)$ in (2.5) so that for all $n \geq n(\epsilon)$ we simultaneously have

$$(2.6) \qquad |n^{-1}\psi_n(\theta^*)| < \min\left\{ -\frac{\log \eta}{2}, \epsilon \right\}$$

and

$$(2.7) \qquad P_n^*\left( |n^{-1}S_{n-k} - \nu| > \frac{\epsilon\nu}{1+\epsilon} \right) \leq \eta^n$$

for $k = 0$ and 1 for some $\eta$ with $0 \leq \eta < 1$. This is possible because of assumption (1.4) and Theorem 7.

For the first term in (2.5),

$$P(T(x) = j) \leq P(S_j > x)$$
$$\leq E_j^*[\exp(-\theta^*S_j + \psi_j(\theta^*)); S_j > x]$$
$$\leq \exp(-\theta^*x)E_j^*[\exp(\psi_j(\theta^*)); S_j > x]$$
$$(2.8) \qquad\qquad\qquad \leq \exp(-\theta^*x)\exp(\psi_j(\theta^*)).$$

We use (1.6) to ensure that (2.8) is finite.

For the second term in (2.5), note that (starting with the reasoning in (2.8))

$$P(T(x) = j) \leq \exp(-\theta^*x)E_j^*[\exp(\psi_j(\theta^*)); S_j > x],$$

where

$$E_j^*[\exp(\phi_j(\theta^*)); S_j > x] \leq \exp\left( j\left( \frac{-\log \eta}{2} \right) \right) P_j^*(S_j > x)$$

$$\leq \exp\left( j\frac{(-\log \eta)}{2} \right) P_j^*\left( S_j > \frac{\nu j}{1 - \epsilon} \right)$$

$$\leq \exp\left( j\frac{(-\log \eta)}{2} \right) P_j^*\left( j^{-1}S_j - \nu > \frac{\nu}{1 - \epsilon} - \nu \right)$$

$$\leq \exp\left( j\frac{(-\log \eta)}{2} \right) P_j^*\left( |j^{-1}S_j - \nu| > \frac{\epsilon\nu}{1 + \epsilon} \right)$$

$$(2.9) \qquad\qquad \leq \exp\left( j\left( \frac{-\log \eta}{2} \right) \right) \eta^j \leq \exp\left( j\frac{\log \eta}{2} \right).$$

Hence,

$$(2.10) \qquad \sum_{j=n(\epsilon)+1}^{\lceil x(1-\epsilon)/\nu \rceil} P(T(x) = j) \leq \exp(-\theta^*x) \sum_{j=1}^{\infty} \eta^{j/2} \leq \exp(-\theta^*x)(1 - \sqrt{\eta})^{-1}.$$

For the third term in (2.5),

$$(2.11) \qquad\qquad P(T(x) = j) \leq \exp(-\theta^*x)E_j^*[\exp \psi_j(\theta^*); S_j > x],$$

where

(2.12)    $E_j^*[\exp(\psi_j(\theta^*)); S_j > x] \le \exp(\psi_j(\theta^*)) \le \exp(j\epsilon) \le \exp(\epsilon x(1+\epsilon)/\nu).$

For the fourth term in (2.5),

$$P(T(x) = j) \le P(S_{j-1} \le x, S_j > x)$$
$$\le E_j^*[\exp(-\theta^* S_j + \psi_j(\theta^*)); S_{j-1} \le x, S_j > x]$$
$$\le \exp(-\theta^* x) \exp(\psi_j(\theta^*)) P_j^*(S_{j-1} \le x),$$

where

$$P_j^*(S_{j-1} \le x) \le P_j\left(S_{j-1} \le \frac{\nu j}{1+\epsilon}\right) = P_j^*\left(\frac{S_{j-1}}{j} \le \frac{\nu}{1+\epsilon}\right)$$
$$\le P_j^*\left(\left|\frac{S_{j-1}}{j} - \nu\right| > \frac{\epsilon\nu}{1+\epsilon}\right) \le \eta^j$$

by Theorem 5 with $k = 1$. Since $\psi_j(\theta^*) \le -(\log\eta)/2$ by (2.6) for $j$ in this sum,

(2.13)    $\displaystyle\sum_{j=\lfloor x(1+\epsilon)/\nu\rfloor}^{\infty} P(T(x) = j) \le \exp(-\theta^* x)\sum_{j=0}^{\infty}\eta^{j/2} \le \exp(-\theta^* x)(1 - \sqrt{\eta})^{-1}.$

Combining (2.5), (2.8), (2.10), (2.12) and (2.13), we obtain

$$P(T(x) < \infty) \le \exp(-\theta^* x)\left\{\sum_{j=1}^{n(\epsilon)}\exp(\psi_j(\theta^*))\right.$$

$$\left. +(1-\sqrt{\eta})^{-1} + \left(\frac{2\epsilon x}{\nu} + 1\right)\exp(\epsilon x(1+\epsilon)/\nu) + (1-\sqrt{\eta})^{-1}\right\}.$$

Hence, using condition (1.6),

$$\limsup_{x\to\infty} x^{-1}\log P(T(x) < \infty) \le -\theta^* + \epsilon(1+\epsilon)/\nu.$$

Since $\epsilon$ was arbitrary,

(2.14)    $\displaystyle\limsup_{x\to\infty} x^{-1}\log P(T(x) < \infty) \le -\theta^*.$

We now establish the lower bound. For this purpose, let $m(\epsilon) = \lceil x(1+\epsilon)/\nu\rceil$. Then

$$P(T(x) < \infty) \ge P(S_{m(\epsilon)} > x)$$
$$\ge E_{m(\epsilon)}^*[\exp(-\theta^* S_{m(\epsilon)} + \psi_{m(\epsilon)}(\theta^*)); S_{m(\epsilon)} > x]$$
$$\ge E_{m(\epsilon)}^*\left[\exp(-\theta^* S_{m(\epsilon)} + \psi_{m(\epsilon)}(\theta^*)); S_{m(\epsilon)} > \frac{m(\epsilon)\nu}{1+\epsilon}\right]$$

$$\geq E^*_{m(\epsilon)}\left[\exp\left(-\theta^* S_{m(\epsilon)} + \psi_{m(\epsilon)}(\theta^*)\right); m(\epsilon)^{-1}|S_{m(\epsilon)} - \nu| < \frac{\epsilon\nu}{1+\epsilon}\right]$$

$$\geq E^*_{m(\epsilon)}\left[\exp\left(-\theta^*\nu\frac{(1+2\epsilon)}{1+\epsilon}m(\epsilon)\right) + \psi_{m(\epsilon)}(\theta^*); |m(\epsilon)^{-1}S_{m(\epsilon)} - \nu| < \frac{\epsilon\nu}{1+\epsilon}\right]$$

$$\geq \exp\left(-\theta^*\nu\frac{(1+2\epsilon)}{1+\epsilon}m(\epsilon) + \psi_{m(\epsilon)}(\theta^*)\right)P^*_{m(\epsilon)}\left(|m(\epsilon)^{-1}S_{m(\epsilon)} - \nu| < \frac{\epsilon\nu}{1+\epsilon}\right).$$

Since

$$P^*_{m(\epsilon)}\left(|m(\epsilon)^{-1}S_{m(\epsilon)} - \nu| < \frac{\epsilon\nu}{1+\epsilon}\right) \to 1 \text{ as } x \to \infty$$

by Theorem 5,

$$\liminf_{x\to\infty} x^{-1}\log P(T(x) < \infty) \geq \liminf_{x\to\infty}\left\{-\theta^*\nu\frac{(1+2\epsilon)}{1+\epsilon}\frac{m(\epsilon)}{x} - \frac{\epsilon m(\epsilon)}{x}\right\}$$

$$\geq -\theta^*(1+2\epsilon) - \epsilon(1+\epsilon)/\nu.$$

Since $\epsilon$ was arbitrary, we conclude that

$$(2.15) \qquad \liminf_{x\to\infty} x^{-1}\log P(T(x) < \infty) \geq -\theta^*.$$

Combining (2.14) and (2.15) completes the proof.

## 3. Proof of Theorem 7

As before, let $\psi_n(\theta) = \log E\exp(\theta S_n)$ and recall that $n^{-1}\psi_n(\theta) \to \psi(\theta)$ as $n \to \infty$ for $|\theta - \theta^*| < \epsilon^*$ where $\theta^*$ and $\epsilon^*$ are as assumed. We start with the case $k = 0$. Then, for each $\theta$, $0 < \theta < \epsilon^*$,

$$P^*_n(n^{-1}S_n > \nu + \epsilon) \leq \exp\left(-\theta n(\nu + \epsilon)\right)E^*_n\exp(\theta S_n),$$

where

$$\exp\left(-\theta n(\nu + \epsilon)\right)E^*_n\exp(\theta S_n) = \exp\left(-\theta n(\nu + \epsilon)\right)$$

$$\times \int \exp(\theta x)\exp(\theta^* x - \psi_n(\theta^*))P(S_n \in dx)$$

$$= \exp\left(-\theta n(\nu + \epsilon)\right)\exp\left(\psi_n(\theta + \theta^*) - \psi_n(\theta^*)\right).$$

We choose $n_0$ in Theorem 7 suitably large so that $\psi_n(\theta + \theta^*) < \infty$, which is possible by (1.4) and (1.5). We use the fact that $E\exp(\theta_1 Z) < \infty$ when $E\exp(\theta_2 Z) < \infty$ and $0 < \theta_1 < \theta_2$, hence

$$(3.1) \qquad \limsup_{n\to\infty} n^{-1}\log P^*_n(n^{-1}S_n > \nu + \epsilon) \leq \psi(\theta^* + \theta) - \psi(\theta^*) - \theta(\nu + \epsilon).$$

However, by Taylor's theorem,

$$\psi(\theta + \theta^*) - \psi(\theta^*) - \theta(\nu + \epsilon) = \psi'(\theta^*)\theta + o(\theta) - \theta(\nu + \epsilon) \text{ as } \theta \to 0$$

$$= -\theta\epsilon + o(\theta) \text{ as } \theta \to 0.$$

Hence, we can choose $\theta_1$ with $0 < \theta_1 < \epsilon^*$ so that

$$\psi(\theta^* + \theta_1) - \psi(\theta^*) - \theta_1(\nu + \epsilon) < -\theta_1\epsilon/2$$

and

$$\limsup_{n \to \infty} n^{-1} \log P_n^*(n^{-1}S_n > \nu + \epsilon) \leq -\theta_1\epsilon/2,$$

which establishes one half of (2.2).

On the other hand, for $0 < \theta < \epsilon^*$,

$$P_n^*(n^{-1}S_n < \nu - \epsilon) = P_n^*(-\theta S_n > -\theta n(\nu - \epsilon))$$

$$\leq \exp(\theta n(\nu - \epsilon))E_n^* \exp(-\theta S_n),$$

where

$$\exp(\theta n(\nu - \epsilon))E_n^* \exp(-\theta S_n) = \exp(\theta n(\nu - \epsilon)) \int \exp(-\theta x) P_n^*(S_n \in dx)$$

$$= \exp(\theta n(\nu - \epsilon)) \exp(\psi_n(\theta^* - \theta) - \psi_n(\theta^*)).$$

Hence,

(3.2) $$\limsup_{n \to \infty} P_n^*(n^{-1}S_n < \nu - \epsilon) \leq \theta(\nu - \epsilon) + \psi(\theta^* - \theta) - \psi(\theta^*).$$

Then, as before,

$$\psi(\theta^* - \theta) - \psi(\theta^*) - \theta(\nu - \epsilon) = -\theta\epsilon + o(\theta) \text{ as } \theta \to 0$$

so that we can choose $\theta_2$ with $0 < \theta_2 < \epsilon^*$ so that

$$\psi(\theta^* - \theta_2) - \psi(\theta^*) + \theta_2(\nu - \epsilon) \leq -\theta_2\epsilon/2$$

and

$$\limsup_{n \to \infty} n^{-1} \log P_n^*(n^{-1}S_n < \nu - \epsilon) \leq -\theta_2\epsilon/2,$$

which completes the proof for $k = 0$.

For $k \geq 1$, we first note that $E \exp(\theta(S_n - S_{n-k})) < \infty$ for all $|\theta| < \delta$ for some $\delta > 0$ if condition (1.8) holds. To see this, apply the Cauchy–Schwarz inequality $k$ times to obtain

$$E \exp(\theta(S_n - S_{n-k})) \leq (E \exp(\theta X_n)^2)^{1/2} E(\exp(\theta(S_{n-1} - S_{n-k}))^2)^{1/2}$$

$$\leq (E \exp(2\theta X_n))^{1/2} E(\exp(2\theta(S_{n-1} - S_{n-k})))^{1/2}$$

$$\leq (E \exp(2\theta X_n))^{1/2} (E \exp(4\theta X_{n-1}))^{1/4} \dots (E \exp(2^{k+1}\theta X_{n-k+1}))^{2^{-(k+1)}}.$$

We choose $n_0$ suitably large so that, for some finite $M$,

$$(E\exp{(2\theta X_n)})^{1/2}E\exp{(4\theta X_{n-1})})^{1/4}\ldots E\exp{(2^{k+1}\theta X_{n-k+1})})^{2^{-(k+1)}} < M$$

for all $n \geq n_0$, which is possible by assumption (1.8).

For $k \geq 1$, we then have

$$\boldsymbol{P}_n^*(n^{-1}S_{n-k} > \nu + \epsilon) \leq \exp{(-\theta n(\nu + \epsilon))}\boldsymbol{E}_n^*\exp{(\theta S_{n-k})}$$

$$\leq \exp{(-\theta n(\nu + \epsilon))}\int \exp\left(\theta \sum_{i=1}^{n-k} x_i\right) \exp\left(\theta^* \sum_{i=1}^{n} x_i - \psi_n(\theta^*)\right)$$

$$\times \boldsymbol{P}(X_1 \in dx_1, \ldots, X_n \in dx_n)$$

$$\leq \exp{(-\theta n(\nu + \epsilon))}\boldsymbol{E}[\exp{((\theta + \theta^*)S_n - \theta(S_n - S_{n-k}) - \psi_n(\theta^*))}]$$

$$\leq \exp{(-\theta n(\nu+\epsilon))}(\boldsymbol{E}\exp{(p(\theta+\theta^*)S_n)})^{1/p}(\boldsymbol{E}\exp{(-q\theta(S_n - S_{n-k})))}^{1/q}\exp{(-\psi_n(\theta^*))}$$

for positive $p$ and $q$ with $p^{-1} + q^{-1} = 1$ by Hölder's inequality. We choose $p$ sufficiently close to 1 and $\theta$ sufficiently small so that $p(\theta + \theta^*)$ is within the required neighborhood of $\theta^*$ and $q\theta < \delta$, so that $\boldsymbol{E}\exp{(-q\theta(S_n - S_{n-k}))}^{1/q}$ is bounded for $n \geq n_0$. Hence,

$$\limsup_{n\to\infty} n^{-1} \log \boldsymbol{P}_n^*(n^{-1}S_{n-k} > \nu + \epsilon) \leq -\theta(\nu + \epsilon) + \psi(p(\theta + \theta^*))^{1/p} - \psi(\theta^*).$$

(3.3)

Since $p$ was arbitrary, we can let $p \to 1$ in (3.3) to obtain the analog of (3.1) with $S_{n-k}$ instead of $S_n$.

Similarly,

$$\boldsymbol{P}_n^*(n^{-1}S_{n-k} < \nu - \epsilon) \leq \exp{(\theta n(\nu - \epsilon))}\boldsymbol{E}_n^*\exp{(-\theta S_{n-k})}$$

$$\leq \exp{(\theta n(\nu - \epsilon))}\int \exp\left(-\theta \sum_{i=1}^{n-k} x_i\right) \exp\left(\theta^* \sum_{i=1}^{n} x_i - \psi_n(\theta^*)\right)$$

$$\times \boldsymbol{P}(X_1 \in dx_1, \ldots, X_n \in dx_n)$$

$$\leq \exp{(\theta n(\nu - \epsilon))}\boldsymbol{E}[\exp{(\theta^* - \theta)S_n + \theta(S_n - S_{n-k}) - \psi_n(\theta^*)}]$$

$$\leq \exp{(\theta n(\nu - \epsilon))}(\boldsymbol{E}\exp{(p(\theta^* - \theta)S_n)})^{1/p}(\boldsymbol{E}\exp{(q\theta(S_n - S_{n-k})))}^{1/q}\exp{(-\psi_n(\theta^*))}$$

for positive $p$ and $q$ with $p^{-1} + q^{-1} = 1$ by Hölder's inequality. Reasoning as in (3.3), we obtain

(3.4) $\displaystyle\limsup_{n\to\infty} n^{-1} \log \boldsymbol{P}_n^*(n^{-1}S_{n-k} < \nu - \epsilon) \leq \theta(\nu - \epsilon) + \psi(p(\theta^* - \theta))^{1/p} - \psi(\theta^*).$

Letting $p \to 1$ in (3.4) we obtain the analog of (3.2) with $S_{n-k}$ instead of $S_n$. The rest of the proof is the same as for $k = 0$.

## 4. Proof of Theorem 4

We construct discrete-time processes satisfying the conditions of Theorem 1 that suitably approximate the continuous-time processes. In particular, for any $\delta > 0$, we construct a discrete-time waiting-time process $\{W_n^\delta\}$ by defining service times $V_n^\delta$ and interarrival times $U_n^\delta$ via

$$(4.1) \qquad V_n^\delta = I((n+1)\delta) - I(n\delta) \quad \text{and} \quad U_n^\delta = \delta, \quad n \geq 1.$$

Since $EY(t) = \rho t$ for $\rho < 1$, $EV_n^\delta < EU_n^\delta$. We initialize by setting $W_0^\delta = L(0)$. Then, by induction, we have

$$(4.2) \qquad W_n^\delta \leq L(n\delta) \leq W_n^\delta + \delta, \qquad n \geq 0.$$

Since $Y(t)$ has stationary increments, $L(t)$ is distributed the same as $\sup\{Y(s) : 0 \leq s \leq t\}$. Since this supremum is non-decreasing, $L(t) \Rightarrow L$ as $t \to \infty$. Since $W_n^\delta \Rightarrow W^\delta$ and $L(n\delta) \Rightarrow L$ as $n \to \infty$, we have

$$(4.3) \qquad E\exp(\theta W^\delta) \leq E\exp(\theta L) \leq \exp(\theta\delta)E\exp(\theta W^\delta).$$

From (4.3), we see that $E\exp(\theta L) < \infty$ if and only if $E\exp(\theta W^\delta) < \infty$.

Hence, it suffices to show that $\sup\{\theta : E\exp(\theta W^\delta) < \infty\} = \theta^*$. For this purpose, let $S_n^\delta = V_0^\delta + \cdots + V_{n-1}^\delta - n\delta$. Then

$$E\exp(\theta(S^{\delta\lfloor t/\delta\rfloor} - \delta)) = E\exp(\theta(Y(\delta\lfloor t/\delta\rfloor) - \delta)) \leq E\exp(\theta Y(t))$$

$$(4.4) \qquad\qquad \leq E\exp(\theta(Y(\delta\lceil t/\delta\rceil) + \delta)) \leq E\exp(\theta(S^{\delta\lceil t/\delta\rceil} + \delta)).$$

Therefore,

$$\limsup_{t\to\infty} t^{-1}\log E\exp(\theta S^{\delta\lfloor t/\delta\rfloor}) \leq \lim_{t\to\infty} t^{-1}\log E\exp(\theta Y(t))$$

$$(4.5) \qquad\qquad\qquad \leq \liminf_{t\to\infty} t^{-1}\log E\exp(\theta S^{\delta\lceil t/\delta\rceil}).$$

Hence,

$$(4.6) \qquad \lim_{n\to\infty} n^{-1}\log E\exp(\theta S_n^\delta) = \delta^{-1}\psi(\theta).$$

Since $\delta^{-1}\psi(\theta) = 0$ if and only if $\psi(\theta) = 0$, the proof is complete.

## 5. Proof of Theorem 5

By Theorem 3,

$$n^{-1}\log E\exp(-\theta S_n^{u*}) \to \psi_u^*(-\theta) = -\psi_A^{-1}(\theta) \text{ as } n \to \infty$$

for each $\theta > 0$. Now we consider the waiting time and workload in the $G/D/1$ queue with the given point process as the arrival process. We let the deterministic service times have mean 1 and the arrival processes have rate $\rho$ where $0 < \rho < 1$. This requires that we scale the original process.

By Theorem 1, in the customer-stationary case (1.1) holds for each $\rho$, $0 < \rho < 1$, where $\theta_W^*(\rho)$ satisfies Equation (1.30), i.e.

$$(5.1) \qquad \psi_u(-\theta_W^*(\rho)/\rho) = -\theta_W^*(\rho), \qquad 0 < \rho < 1.$$

On the other hand, by Theorem 4 and Proposition 7, in the time-stationary case (1.36) holds for each $\rho$, $0 < \rho < 1$, where the decay rate $\theta_L^*(\rho)$ satisfies

$$(5.2) \qquad \rho\psi_A(\psi_v(\theta_L^*(\rho))) - \theta_L^*(\rho) = \rho\psi_A(\theta_L^*(\rho)) - \theta_L^*(\rho) = 0,$$

because the decay rate function of $A(\rho t)$ is $\rho\psi_A$. However, by (1.29), (5.2) is equivalent to

$$(5.3) \qquad -\psi_u^*(-\theta_L^*(\rho)/\rho) = \psi_A^{-1}(\theta_L^*(\rho)/\rho) = \theta_L^*(\rho), \qquad 0 < \rho < 1.$$

By Proposition 8, $\theta_W^*(\rho) = \theta_L^*(\rho)$ for $0 < \rho < 1$. Hence, (5.3) becomes

$$(5.4) \qquad \psi_u^*(-\theta_W^*(\rho)/\rho) = -\theta_W^*(\rho), \qquad 0 < \rho < 1.$$

Finally, by Proposition 5, (5.1) and (5.4) imply that $\psi_u = \psi_u^*$.

## 6. Proof of Theorem 6

By Theorem 5, $\psi_u(-\theta) = -\psi_A^{-1}(\theta)$. Paralleling the proof of Theorem 5, we have

$$(6.1) \qquad \psi_u(-\theta_W^*(\rho)/\rho) = -\psi_v(\theta_W^*(\rho)), \qquad 0 < \rho < 1,$$

instead of (5.1) and

$$(6.2) \qquad \rho\psi_A(\psi_v(\theta_L^*(\rho)) - \theta_L^*(\rho) = 0, \qquad 0 < \rho < 1,$$

instead of (5.2). However, (6.2) is equivalent to

$$(6.3) \qquad \psi_A^{-1}(\theta_L^*(\rho)/\rho) = \psi_v(\theta_L^*(\rho)), \qquad 0 < \rho < 1.$$

Since $\psi_A^{-1}(\theta) = -\psi_u(-\theta)$, (6.3) coincides with (6.1), so that we must have $\theta_L^*(\rho) = \theta_W^*(\rho)$, $0 < \rho < 1$.

## 7. Proof of Proposition 7

Note that, for any $\epsilon > 0$, there is an $n_0$ such that

$$E \exp(\hat{\theta}I(t)) = \sum_{n=0}^{\infty} E \exp\left(\hat{\theta}\sum_{i=1}^{n} V_i\right) P(A(t) = n)$$

$$\leq \sum_{n=0}^{\infty} \exp(n(\psi_v(\hat{\theta}) + \epsilon)) P(A(t) = n) + E \exp\left(\hat{\theta}\sum_{i=0}^{n_0} V_i\right)$$

$$\leq E \exp\left(\psi_v(\hat{\theta}) + \epsilon\right) A(t) + E \exp\left(\hat{\theta} \sum_{i=0}^{n_0} V_i\right)$$

$$\leq \exp\left(t\psi_A(\psi_v(\hat{\theta}) + \epsilon) + \epsilon\right) + E \exp\left(\hat{\theta} \sum_{i=0}^{n_0} V_i\right)$$

for $t$ suitably large. Hence,

$$\limsup_{t \to \infty} t^{-1} \log E \exp\left(\hat{\theta} I(t)\right) \leq \psi_A(\psi_V(\hat{\theta}) + \epsilon) + \epsilon.$$

Since $\epsilon$ was arbitrary and $\psi_A$ is continuous at $\psi_v(\hat{\theta})$,

$$\limsup_{t \to \infty} t^{-1} \log E \exp\left(\hat{\theta} I(t)\right) \leq \psi_A(\psi_V(\hat{\theta})).$$

The reasoning for the other direction is essentially the same.

## 8. Proof of Proposition 9

We shall work with characterization (ii) in Proposition 1. Note that

$$(8.1) \qquad W = \sum_{i=1}^{(Q^a - 1)^+} V_i + V_e^a \geq \sum_{i=1}^{(Q^a - 1)^+} V_i$$

and

$$(8.2) \qquad L = \sum_{i=1}^{(Q - 1)^+} V_i + V_e \geq \sum_{i=1}^{(Q - 1)^+} V_i,$$

where $V_e^a$ and $V_e$ are the equilibrium residual service times of the customer in service (which in general depend upon $Q^a$ and $Q$). Since the argument is essentially the same for $W$ and $L$, we henceforth consider only $W$. To have a useful inequality in the opposite direction, we truncate the service times by setting $V_n^c = \min\{V_n, c\}$, $n \geq 1$. Then

$$(8.3) \qquad W^c \leq \sum_{i=1}^{Q^{ac}} V_i^c + c.$$

From (8.1), we obtain

$$E \exp(\theta W) \geq E(E \exp(\theta V_1))^{Q^a - 1},$$

so that $\theta_{Q^a}^* \geq \log E \exp(\theta_W^* V_1)$. (As in Proposition 8, we use the fact that $\theta_W^* < \theta_{V_1}^*$.) From (8.3), we obtain

$$E \exp(\theta W^c) \leq \exp(\theta_c) E(E \exp(\theta V_1^c))^{Q^{ac}},$$

so that

$$\log E \exp(\theta_{W^c}^* V_1^c) \geq \theta_{Q^{ac}}^*.$$

Then note that $V_1^c$, $W^c$, $Q^{ac}$, $Q^c$ and $L^c$ increase stochastically to their limits $V_1$, $W$, $Q^a$, $Q$ and $L$ as $c \to \infty$, see Theorems 4, 5 and 8 plus the remark on p. 216 of Whitt (1981). Hence

$$\theta_{Q^a}^* \leq \liminf_{c \to \infty} \theta_{Q^c}^* \leq \lim_{c \to \infty} \log E \exp(\theta_{W^c}^* V_1^c) = \log E \exp(\theta_W^* V_1) < \infty.$$

## 9. An example

In this section we give an example.

*Example* 1. To see that the conditions in Theorem 1 are not necessary for (1.1) or (1.2), consider the $G/G/1$ model with

$$P(U_{2n+1} = 1, U_{2n+2} = 1 + Y_n, V_{2n+1} = 1 + Y_n, V_{2n+2} = 0 \text{ for all } n) = 1/2$$

and

$$P(U_{2n+1} = 1 + Y_n, U_{2n+2} = 1, V_{2n+1} = 0,$$

$$V_{2n+2} = 1 + Y_{n+1} \text{ for all } n) = 1/2,$$

where $\{Y_n\}$ is an i.i.d. sequence of exponential random variables with mean 1. Then $\{U_n, V_n\}$ is stationary with $EV_n = 1 < EU_n = 3/2$, so that $\rho = 2/3$. Moreover, it is easy to see that, for $n \geq 1$,

$$P(W_n > x) = P(W > x) = (1/2)e^{-x}, \qquad x > 0,$$

but

$$P(S_{2n+1} = Y_n, S_{2n+2} = -n \text{ for all } n)$$

$$= P(S_1 = -(n + Y_1), S_{2n+2} = Y_{n+1} - (n + Y_1) \text{ for all } n) = 1/2,$$

so that

$$E \exp(\theta S_{2n+1}) = \tfrac{1}{2} E \exp(\theta Y_1) + \tfrac{1}{2} E \exp(-\theta(n + Y_1)),$$

$$E \exp(\theta S_{2n+2}) = \tfrac{1}{2} \exp(-\theta n) + \tfrac{1}{2} E \exp(\theta(Y_{n+1} - n - Y_1)),$$

and

$$n^{-1} \log E \exp(\theta S_n) \to \psi(\theta) = -\theta/2 \text{ as } n \to \infty.$$

Hence, (1.1) and (1.2) hold, but $\psi(\theta^*) = 0$ for $\theta^* = 0$.

*Note added in proof*: Additional results for multiserver queues related to those mentioned on p. 132 appear in Sadowsky and Szpankowski (1995) and Sadowski (1995).

## References

ABATE, J. AND WHITT, W. (1994) A heavy-traffic expansion for asymptotic decay rates of tail probabilities in multi-channel queues. *Operat. Res. Lett.* To appear.

ABATE, J., CHOUDHURY, G. L. AND WHITT, W. (1994a) Exponential approximations for tail probabilities of queues, I: waiting times. *Operat. Res.* To appear.

ABATE, J., CHOUDHURY, G. L. AND WHITT, W. (1994b) Exponential approximations for tail probabilities of queues, II: sojourn time and workload. *Operat. Res.* To appear.

ABATE, J., CHOUDHURY, G. L. AND WHITT, W. (1994c) Asymptotics for steady-state tail probabilities in structured Markov queueing models. *Stoch. Models* **10**. To appear.

ABATE, J., CHOUDHURY, G. L. AND WHITT, W. (1993) Calculation of the $GI/G/1$ waiting-time distribution and its cumulants from Pollaczek's formulas. *AEÜ* **47**, 311–321.

ADDIE, R. G. AND ZUCKERMAN, M. (1993) An approximation for performance evaluation of stationary single server queues. *IEEE Trans. Commun.*, to appear.

ASMUSSEN, S. (1987) *Applied Probability and Queues.* Wiley, New York.

ASMUSSEN, S. (1989) Risk theory in a Markovian environment. *Scand. Actuarial J.* 69–100.

ASMUSSEN, S. (1993) *Fundamentals of Ruin Probability Theory.* Preliminary draft, Aalborg University, Denmark.

ASMUSSEN, S. AND PERRY, D. (1992) On cycle maxima, first passage problems and extreme value theory for queues. *Stoch. Models* **8**, 421–458.

BOROVKOV, A. A. (1976) *Stochastic Processes in Queueing Theory.* Springer-Verlag, New York.

BUCKLEW, J. A. (1990) *Large Deviations Techniques in Decision, Simulation and Estimation.* Wiley, New York.

CHANG, C.-S. (1992) Stability, queue length and delay, part II: stochastic queueing networks. *Proc. 31st IEEE Conf. Decision and Control*, Tucson, AZ, 1005–1010.

CHANG, C.-S. (1993) Stability, queue length and delay of deterministic and stochastic queueing networks. IBM T. J. Watson Research Center, Yorktown Heights, NY.

CHANG, C.-S., HEIDELBERGER, P., JUNEJA, S. AND SHAHABUDDIN, P. (1992) Effective bandwidth and fast simulation of ATM intree networks. IBM T. J. Watson Research Center, Yorktown Heights, NY.

CHOUDHURY, G. L. AND WHITT, W. (1994) Heavy-traffic asymptotic expansions for the asymptotic decay rates in the $BMAP/G/1$ queue. *Stoch. Models.* To appear.

CHOUDHURY, G. L., LUCANTONI, D. M. AND WHITT, W. (1993) Squeezing the most out of ATM. Submitted.

CHOUDHURY, G. L., LUCANTONI, D. M. AND WHITT, W. (1994) On the effectiveness of effective bandwidths for admission control in ATM networks. *Proc. 14th Internat. Teletraffic Congr.*, Antibes, France. To appear.

CHUNG, K. L. (1974) *A Course in Probability Theory*, 2nd edn. Academic Press, New York.

DEMBO, A. AND ZEITOUNI, O. (1992) *Large Deviations Techniques and Applications.* A. K. Peters, Wellesley, MA.

ELLIS, R. (1984) Large deviations for a general class of random vectors. *Ann. Prob.* **12**, 1–12.

ELWALID, A. I. AND MITRA, D. (1993) Effective bandwidth of general Markovian traffic sources and admission control of high speed networks. *IEEE/ACM Trans. Networking* **1**, 329–343.

ELWALID, A. I. AND MITRA, D. (1994) Markovian arrival and service communication systems: spectral expansions, separability and Kronecker-product forms. AT&T Bell Laboratories, Murray Hill, NJ.

FELLER, W. (1971) *An Introduction to Probability Theory and its Applications*, Vol. II, 2nd edn. Wiley, New York.

FRANKEN, P., KÖNIG, D., ARNDT, U. AND SCHMIDT, V. (1981) *Queues and Point Processes.* Akademie-Verlag, Berlin.

GÄRTNER, J. (1977) On large deviations from the invariant measure. *Theory Prob. Appl.* **22**, 24–39.

GIBBENS, R. J. AND HUNT, P. J. (1991) Effective bandwidths for the multi-type UAS channel. *QUESTA* **9**, 17–28.

GLYNN, P. W. AND WHITT, W. (1988a) Ordinary CLT and WLLN versions of $L = \lambda W$. *Math. Operat. Res.* **13**, 674–692.

GLYNN, P. W. AND WHITT, W. (1988b) An LIL version of $L = \lambda W$. *Math. Operat. Res.* **13**, 693–710.

GLYNN, P. W. AND WHITT, W. (1994) Large deviations behavior of counting processes and their inverses. *QUESTA*. To appear.

GUERIN, R., AHMADI, H. AND NAGHSHINEH, M. (1991) Equivalent capacity and its application to bandwidth allocation in high-speed networks. *IEEE J. Sel. Areas Commun.* **9**, 968–981.

IGLEHART, D. L. AND WHITT, W. (1970) Multiple channel queues in heavy-traffic, II: sequences, networks and batches. *Adv. Appl. Prob.* **2**, 355–369.

IGLEHART, D. L. AND WHITT, W. (1971) The equivalence of functional central limit theorems for counting processes and associated partial sums. *Ann. Math. Statist.* **42**, 1372–1378.

KEILSON, J. (1965) *Greens' Function Methods in Probability Theory.* Griffin, London.

KELLY, F. P. (1991) Effective bandwidths at multi-class queues. *QUESTA* **9**, 5–16.

KINGMAN, J. F. C. (1962) On queues in heavy traffic. *J. R. Statist. Soc.* B **24**, 383–392.

NEUTS, M. F. (1986) The caudal characteristic curve of queues. *Adv. Appl. Prob.* **18**, 221–254.

SHWARTZ, A. AND WEISS, A. (1993) *Large Deviations for Performance Analysis: Queues, Communication and Computing.* In preparation.

SMITH, W. L. (1953) On the distribution of queueing times. *Proc. Camb. Phil. Soc.* **49**, 449–461.

SOHRABY, K. (1993) On the theory of general on-off sources with applications in high speed networks. Proc. *INFOCOM '93*, San Francisco. pp. 401–410.

STERN, T. E. AND ELWALID, A. I. (1991) Analysis of a separable Markov-modulated rate model for information-handling systems. *Adv. Appl. Prob.* **23**, 105–139.

TAKÁCS, L. (1962) *Introduction to the Theory of Queues.* Oxford University Press, New York.

TAKÁCS, L. (1963) The limiting distribution of the virtual waiting time and the queue size for a single server queue with recurrent input and general service times. *Sankhyā* A **25**, 91–100.

TAKÁCS, L. (1967) *Combinatorial Methods in the Theory of Stochastic Processes.* Wiley, New York.

TIJMS, H. C. (1986) *Stochastic Modeling and Analysis: A Computational Approach.* Wiley, New York.

VAN OMMEREN, J. C. W. (1988) Exponential expansion for the tail of the waiting-time probability in the single-server queue with batch arrivals. *Adv. Appl. Prob.* **20**, 880–895.

WHITT, W. (1980) Some useful functions for functional limit theorems. *Math. Operat. Res.* **5**, 67–85.

WHITT, W. (1981) Comparing counting processes and queues. *Adv. Appl. Prob.* **13**, 207–220.

WHITT, W. (1994) Tail probabilities with statistical multiplexing and effective bandwidths for multi-class queues. *Telecommunication Systems*, to appear.

*References added in proof*

SADOWSKY, J. S. (1995) The probability of large queue lengths and waiting times in a heterogeneous multiserver queue, part II: positive recurrence and logarithmic limits. *Adv. Appl. Prob.* **27**(2). To appear.

SADOWSKY, J. S. AND SZPANKOWSKI, W. (1995) The probability of large queue lengths and waiting times in a heterogeneous multiserver queue, part I: tight limits. *Adv. Appl. Prob.* **27**(2). To appear.