# The Asymptotic Behavior of a Network Multiplexer with Multiple Time Scale and Subexponential Arrivals

Predrag R. Jelenković and Aurel A. Lazar

ABSTRACT  Real-time traffic processes, such as video, exhibit multiple time scale characteristics, as well as subexponential first and second order statistics. We present recent results on evaluating the asymptotic behavior of a network multiplexer that is loaded with such processes.

## 1  Introduction

One of the key features in Asynchronous Transfer Mode (ATM) based broadband networks is statistical multiplexing (SMUX). Most of the multiplexed entities are calls originating from various sources. In order to operate properly, each of these calls has to satisfy some quality of service requirements (QOS). QOS requirements are usually bounds on performance measures characterizing the dynamic behavior of the multiplexed traffic. The most basic model of a SMUX is an infinite buffer single server queue with a work conserving scheduler. The fundamental performance measure is the queue length distribution ($\mathbb{P}[Q > x]$). Therefore, it is of utmost importance to have feasible procedures for calculating this distribution under reasonable assumptions on the arrival processes.

Numerous investigations have shown that the arrival processes (sources) that arise in ATM networks (like voice and video) have a very complex statistical structure; an especially troublesome characteristic is the high statistical dependency (e.g., see [25, 30]). Modeling of this high dependency usually leads to analytically very complex statistical characteristics, typically making the associated evaluation of the queue length distribution intractable. However, because of the stringent QOS requirements in ATM, only the tail of the queue length distribution in the domain of very

small probabilities is needed. This has motivated researchers to investigate possible approximations of the asymptotic behavior of the queue length distribution. This is the main subject of our presentation.

More formally, given an infinite buffer single-server queue, let $A = \{A_t, t \geq 0\}$, $C = \{C_t, t \geq 0\}$, be two discrete time, stationary, and ergodic processes (on a probability space $(\Omega, \mathcal{F}, \mathbb{P})$); $A_t$ represents the amount of arrivals to the queue at time $t$, and $C_t$ is the server capacity at time $t$. Then, for any initial random variable $Q_0$, the following (Lindley's) equation

$$Q_{t+1} = (Q_t + A_t - C_t)^+ \tag{1.1}$$

completely defines the queue length process $\{Q_t, t \geq 0\}$. Queues of this type represent a natural model for ATM multiplexers. According to the classical result of Loynes [31], if $\mathbb{E}A_t < \mathbb{E}C_t$ (and $\{A_t, C_t, t \geq 0\}$ are stationary and ergodic) $\{Q_t\}$ couples with the unique stationary solution $\{Q_t^s\}$ of the recursion (1.1) for any initial condition $Q_0$; in particular $\mathbb{P}[Q_t \geq x] \rightarrow \mathbb{P}[Q_0^s \geq x]$ as $t \rightarrow \infty$ (for simplicity we will refer to $Q_t^s$ simply as $Q$). In what follows the difference between the arrival and the service process $\{X_t \overset{def}{=} A_t - C_t, t \geq 0\}$ will be called the *queue increment process*.

Stationarity and ergodicity comprise the general framework for our current exposition, and will be assumed in the rest of the paper. We will see that under different assumptions on the distribution of the queue increment process, the queue length asymptotics may exibit very different behavior. Two major probabilistic categories of assumptions are the *exponential* (Cramér) and *subexponential*. Informally, the exponential category is represented with random variables whose moment generating functions are finite in some positive neighborhood of zero, whereas the subexponential category consists of random variables whose m.g.f.s are infinite on the positive real axis. This paper is organized according to this categorization.

In the first part of the paper (section 2) we examine the exponential asymptotic queueing behavior in the presence of multiple time scales. We demonstrate that in this case the dominant (or so called Equivalent Bandwith) multiplexer approximation may be very inaccurate. To try to alleviate this problem, in section 2.1 we present an asymptotic expansion approach for approximating all queue length probabilities for the case of structured Markovian multiple time scale (decomposable) arrivals. In section 2.2 we prove that for arrival processes that spend long-tailed (random) time in their high activity states, the Equivalent Bandwith (EB) constant does not depend on slow time scale statistics and is equal to the case when processes stay in high activity states all of the time.

In the second part of the paper (section 3) we discuss the problem of approximating the queue length probabilities under subexponential (non Cramér) assumptions. We first give precise definitions and some intuition behind the modeling of real time processes using subexponential statistics. Some very recent asymptotic results for arrival processes with both

subexponential marginals and subexponetial autocorrelation function are summarized in section 3.1. The paper is concluded in section 4.

## 2  Multiple Time Scale Arrivals

Very often, arrival processes that arise in modern communication networks exhibit a multiple time scale structure. A typical example is Variable Bit Rate (VBR) video traffic. This traffic consists of ATM cells, that, when grouped together, correspond to slices; slices are the building blocks for frames, and finally, a large number of frames form scenes [30]. Each of these VBR video building blocks (cells, slices, frames, scenes) belong to different time scales, and are characterized by different statistics. Furthermore, on an even larger time scale these building blocks form calls with their own statistics. The call statistics themselves may change according to the time of the day. Thus, from this brief analysis, we see that there is a wide spectrum of time scales that are involved in modeling flows in broadband networks. The total range is from a few nanoseconds ($ns$) to a few hours ($1hour = 3.6 \ 10^{12} ns$). In this section we will attempt to answer some questions on the queue length asymptotics in the presence of multiple time scale arrivals and Gärtner-Ellis (Cramér) assumptions.

Using the Theory of Large Deviations (see [37]), under general assumptions of the Gärtner-Ellis (Cramér) type, one can show that

$$\lim_{x \to \infty} -\frac{\log \mathbb{P}[Q > x]}{x} = \theta^*, \tag{1.2}$$

for some positive constant $\theta^*$, called the *asymptotic decay rate* (or the equivalent bandwidth constant) [7, 17]. Also, in some cases, like finite Markov arrival and service processes, the following stronger result holds: $\mathbb{P}[Q > x] \sim \alpha e^{-\theta^* x}$ as $x \to \infty$, where $\alpha$ and $\theta^*$ are positive constants; $\theta^*$ is the same as in (1.2). For simple arrival processes (like On-Off Markov sources) it turns out that the constant $\alpha$ is of the order one. This led many authors to believe that the simple approximation $\mathbb{P}[Q > x] \approx e^{-\theta^* x}$ holds; this approximation is commonly referred to as [9] the effective bandwidth (EB) approximation (sometimes it is also called the dominant root approximation). Following this result admission control policies based on the concept of effective bandwidth have been developed; see [7, 16, 18, 17, 27].

However, as discussed in [9], the EB approximation may often be very inaccurate. This is usually the case when many sources ($N$) are multiplexed; under this assumption it was shown in [9] that $\alpha \approx e^{-\gamma N}$ for some constant $\gamma$. A more formal analysis of the multiplexing of a large number of sources and an improvement of the EB approximation is given in [14, 15]. Complementing the work done in [9], in [19] we have shown that EB approximation may be very inaccurate in the presence of multiple time scale

arrivals. Similar observations of inaccuracy of the EB approximation in the presence of multiple time scales (in the context of nearly decomposable Markov-modulated arrivals) were independently obtained in [35].

*From a mathematical point of view, the inaccuracy of the EB approximation is due to fact that two processes that are "close" in the distribution sense may be far apart in the cumulant sense.* Recall that a family of processes is said to converge in distribution if *all finite dimensional* distributions of these processes converge in distribution. On the other hand, the asymptotic decay rate constant $\theta^*$ is completely determined by the cumulant function $\varphi(\theta)$ which is a functional of the *whole (infinite dimensional)* arrival process. Therefore, convergence in distribution of a family of processes does not necessarily imply convergence of their cumulant functions.

A simple numerical example with two state Markov-modulated arrivals that illustrate the preceding comments (on the disagreement of the two convergence concepts) is given in [19]. The structure of the example is as follows. The modulating chain is assumed to have two states (say $\{1, 2\}$) with transition probabilities $p_{21} = \epsilon, p_{12} = o(\epsilon)$; when in state $i = 1, 2$ the source is producing i.i.d. arrivals $Y(j)$ such that in state 1 the source is producing stochastically smaller arrivals than in state 2. Since, $p_{12} = o(\epsilon)$ it is easy to see that the arrival process converges in distribution to the stochastically smaller process $Y(1)$, as $\epsilon \to 0$. However, its cumulant function converges to the cumulant function of the stochastically larger process $Y(2)$ (a formal argument that justifies this can be found in the proof of Theorem 2.1). For different values of the paramenter $\epsilon$ queue probabilities are presented in Figure 1 (solid lines). We can see that as $\epsilon \to 0$ the queue distribution converges to the queue distribution when $A \equiv Y(1)$ (represented by the common steep decline of the three solid lines on Figure 1), but the tail always decays as if the arrival process is the stochastically larger process $A \equiv Y(2)$ (parallel lines). Also note that the EB approximation (dashed line on the same figure) is off by orders of magnitude from the true probabilities. (For more details and more examples see [19].)

This idea was exploited in greater generality in [19], where, for a family of arrival processes $A^\epsilon, \epsilon > 0$, we give sufficient conditions under which the queue length distribution satisfies the following extension of the logarithmic asymptotic relation (1.2)

$$\lim_{\epsilon \to 0} \overline{\lim}_{r \to \infty} \frac{-\log \mathbb{P}[Q^\epsilon \geq x]}{x} = \theta^*,$$

for some $\theta^* > 0$; symbol $\overline{\lim}$ denotes that either lim sup or lim inf is taken. Using this result we have shown, under strict stability conditions, that the asymptotic decay rate of an ATM multiplexer *does not depend on the slow time scale statistics* (larger time units). However, the rate at which the queue length distribution decreases for small buffer sizes could be much larger than the asymptotic decay rate. This implies that an equivalent bandwidth admission control policy (based only on $\theta^*$) may *significantly*

*underutilize the system resources*, and that the slower timescales can be very important here.
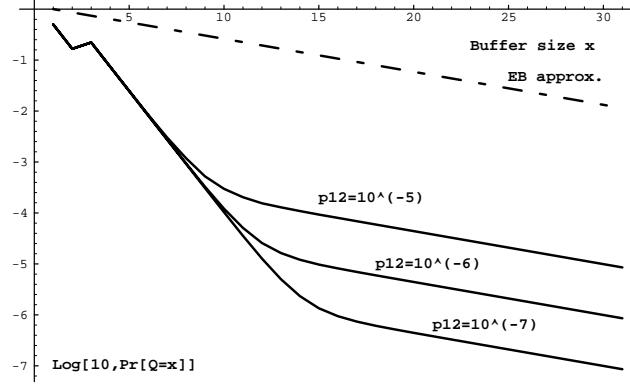


FIGURE 1. Graph of $\log_{10} \mathbb{P}[Q^\epsilon = x]$ from example 1.

In the same paper it was experimentally confirmed that the histogram of the queue distribution obtained by statistical multiplexing of 6 parts of the Star Wars video sequence has a *"polygonal shape"* (multiple decay rates), typical for multiple time scale models, see Figure 2.
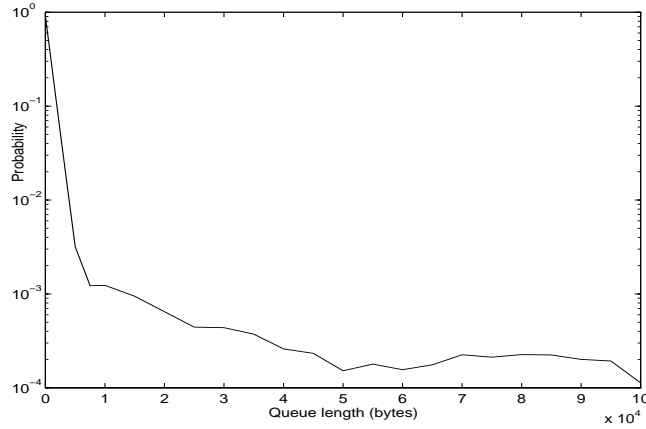


FIGURE 2. Queue length distribution for multiplexing 6 parts of the Star Wars video sequence on the slice level. The total length of the multiplexed sequence is 1,000,000 slices ($\approx$ 23 min).

## 2.1   Queueing Analysis

Overall, as shown in [19], the EB approximation may give very inaccurate results. Also, as pointed out in [9], the exact asymptotic single exponential approximation may be poor (the authors suggested a procedure for approximating the queue length distribution with three exponentials). For that reason we have investigated a perturbation theory based approach for approximating all queue probabilities in the presence of multiple time scale (nearly decomposable) [20] arrival processes (a comprehensive treatment of a discrete time queue with multiple time scale arrivals can be found in [21]).

In that work we developed a *recursive* asymptotic expansion method for approximating the queue length distribution and investigated the radius of convergence of the queue asymptotic expansion series. The analysis focused on "small" to "moderate" buffer sizes under the conditions of strictly stable multiple time scale arrivals. For a class of examples we *analytically* determined the radius of convergence using methods of linear operator theory. We also gave general sufficient conditions under which the radius converges to zero; this showed roughly what situations have to be avoided for the proposed method to work (well). We combined the asymptotic expansion method with the EB approximation, and gave an approximation procedure for the buffer probabilities for all buffer ranges. The procedure was tested on extensive numerical examples. We illustrate this procedure in the following numerical example.

The asymptotic expansion approximation with $k = 0, 3$ expansion terms for multiplexing 8 heterogeneous On-Off sources is displayed in Figure 3 (capacity of the server was taken to be $C_t \equiv 1$). The combination of asymptotic expansion and EB approximation is plotted in Figure 4. We see that the transition between the two approximations is smooth. Therefore, although we have no error estimate in the EB domain, from the smoothness of transition, we can expect that the approximation is excellent in the EB domain as well. This *smoothness of fit* can be used as a heuristic criterion for the overall accurateness of the approximation.

Let us now compare this approximative method with the classical exact z-transform inversion. In order to obtain the exact solution, one must find the inverse of $64X64$ z-transform matrix, then find 63 roots of the characteristic polynomial in the unit circle; use these roots to obtain boundary probabilities, and, at last, find the inverse z-transform of the queue z-transform. We were not able to complete even the first step, i.e., finding the z-transform matrix inverse after 24 hours, after which we stopped the program. (Computation was attempted with Mathematica 2.2 on a (150MHz, 64M RAM, 100M virtual memory)) SGI machine. However, using the same environment Mathematica + SGI, we obtained a three term expansion approximation in less than an hour. This clearly shows the efficacy of the asymptotic expansion method.
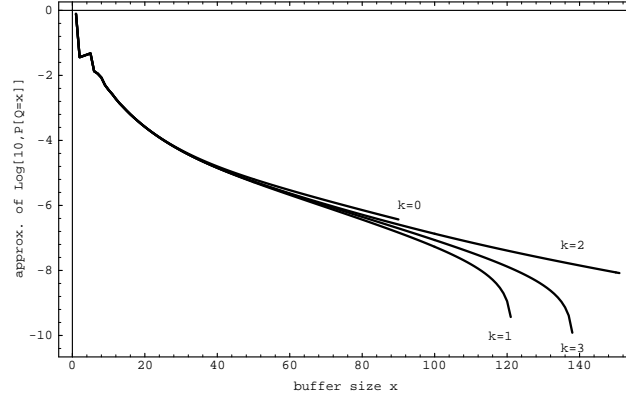
FIGURE 3. Approximate "probabilities" obtained by using $k$ expansion terms.

As all of the calculations were done with Mathematica 2.2, which is known to be slow for intensive numerical problems, we expect that the asymptotic expansion method when implemented in C will produce much faster results. Therefore, we predict that this method will be very useful for large practical problems that often appear in the fine tuning of ATM admission controllers.
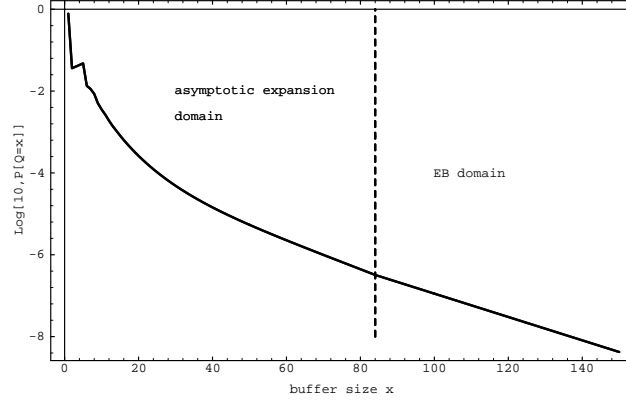


FIGURE 4. Total queue distribution approximation, obtained by combination of the asymptotic expansion method and EB approximation.

We next present a class of arrivals that also cannot be well approximated with an EB approximation. Unlike the case of nearly decomposable Markov arrivals, for which the asymptotic expansion method for approximating queue probabilities is available, approximating the queue probabilities for this type of arrivals remains an open problem.

## 2.2   Arrivals with Long-tailed High Activity Periods

In this section we will examine arrival processes that stay in their high activity states for a long-tailed (random) period of time. For these arrivals we prove that the EB constant does not depend on the modulating process statistics, and is the same as in the worst case when arrival processes are in their high activity states all the time. As a device for modeling time scales, we will consider modulating arrivals of the form $A_t = Y_t(B_t)$, where intuitively $Y_t$ can be thought of as representing the fast time scale changes in the arrival process, and $B_t$ as representing the slow time scale changes.

More formally, let there be $N$ traffic sources all modeled as stationary, ergodic, discrete time stochastic processes $\{A_t^i, t \geq 0, 1 \leq i \leq N\}$. For each $i, 1 \leq i \leq N$, we define

$$A_t^i \stackrel{def}{=} Y_t^i(B_t^i),$$

where $\{Y^i(1), \ldots, Y^i(K_i)\}, K_i \geq 1$ are stationary ergodic processes that are stochastically ordered such that $Y^i(j) \leq_{st} Y^i(K_i), 1 \leq j < K_i$ (for stochastic ordering see [34], or [4], chapter 4). Further, the processes $B^i = \{B_t^i, t \geq 0\}$ are stationary, ergodic, discrete time process with a finite state space $\mathcal{S}_i = \{1, \ldots, K_i\}$. All processes $Y^i(j)$, and $B^i$ are assumed independent of each other.

We assume that each modulating process $B^i$, once in its largest state $K_i$, stays there for a random amount of time with a long-tailed distribution (see Definition 1.8 in the following section). Examples of long-tailed distributions are Pareto, some Weibull, and lognormal; for more examples see the following section. Note that long-tailed distributions decay more slowly than any exponential; in particular, the moment generating function of a long-tailed distribution is infinite on the positive real axis.

In order to state our result we need the following Large Deviations Theory definitions. For $\theta \in \mathbb{R}^+$ let us define the *cumulant function*

$$\varphi_n(\theta) \equiv \varphi_n^A(\theta) \stackrel{def}{=} \frac{1}{n} \log \mathbb{E}\{\exp[\theta \sum_{t=1}^{n} A_t]\}, \tag{1.3}$$

where $n \geq 1$, and let

$$\varphi(\theta) \stackrel{def}{=} \lim_{n \to \infty} \varphi_n(\theta). \tag{1.4}$$

Furthermore, let us define

$$D \stackrel{def}{=} \{\theta \geq 0 \ : \ \varphi(\theta) < \infty\},$$

and make usual Large Deviation **Assumptions:**

**A1** $\varphi(\theta)$ is strictly convex on $D$,

**A2** $\varphi(\theta)$ is differentiable for all $\theta \in D$.

For the sake of simplicity, in this section we will assume that the server capacity process $C_t \equiv c \in \mathbb{R}^+$. This is frequently the case in communication networks.

**Theorem 2.1** Assume that for each $i$ and $\theta \in \mathbb{R}^+, \varphi^{Y^i(K_i)}(\theta)$ exist, and satisfies conditions **A1** and **A2** with $\sum_{i=1}^{N}(\varphi^{Y^i(K_i)})'(0) < c$. Also, let $\pi_{K_i} = \mathbb{P}[B_t^i = K_i] > 0$ for all $i$, and assume that the residual time of staying in the high activity state $K_i$ is long-tailed. Then

$$\lim_{x \to \infty} -\frac{\log \mathbb{P}[Q > x]}{x} = \theta_N^*, \tag{1.5}$$

where $\theta_N^*$ is the equivalent bandwidth constant. This constant, if it exists, is the positive solution of the equation $\sum_{i=1}^{N} \varphi^{Y^i(K_i)}(\theta) = \theta c$; if the positive solution of this equation does not exists, we set $\theta_N^* = \infty$.

*Proof.* Let $T_{K_i}$ be the residual time of staying in state $K_i$; From the assumption that $T_{K_i}$ is long-tailed and Lemma 3.4 it follows that

$$\lim_{n \to \infty} \frac{\log \mathbb{P}[T_{K_i} > n]}{n} = 0. \tag{1.6}$$

Now the theorem will follow from Theorem 3.9 in [7] if we prove that the cumulant function of the arrival process $i$ is equal to the cumulant function of the process $Y^i(K_i)$, i.e., when the arrival process is in its highest activity period $K_i$ all the time.

First, from the stochastic ordering of $Y^i(k), 1 \le k \le K_i$, and Strassen's theorem (see section 4.2.3 [4]) it follows that $Y^i(k)$ can be constructed on the same probability space, such that $Y^i(j) \le Y^i(K_i), 1 \le j < K_i$, holds along each sample path. From this, it follows that

$$\overline{\varphi}^{A^i}(\theta) \le \varphi^{Y^i(K_i)}(\theta). \tag{1.7}$$

The lower bound follows from

$$\liminf_{T \to \infty} \frac{1}{T} \log \mathbb{E}\left[e^{\theta \sum_{t=1}^{T} A_t^i}\right]$$

$$\ge \liminf_{T \to \infty} \frac{1}{T} \log \left(\pi_{K_i}^i \mathbb{P}[T_{K_i} > T]\mathbb{E}\left[e^{\theta \sum_{t=1}^{T} Y_t^i(K_i)}\right]\right)$$

$$= \liminf_{T \to \infty} \frac{\log \mathbb{P}[T_{K_i} > T]}{T} + \varphi^{Y^i(K_i)}(\theta) = \varphi^{Y^i(K_i)}(\theta),$$

where the last equality follows from (1.6). This proves that $\varphi^{A^i}(\theta) = \varphi^{Y^i(K_i)}(\theta)$, and the assertion of the theorem follows.□

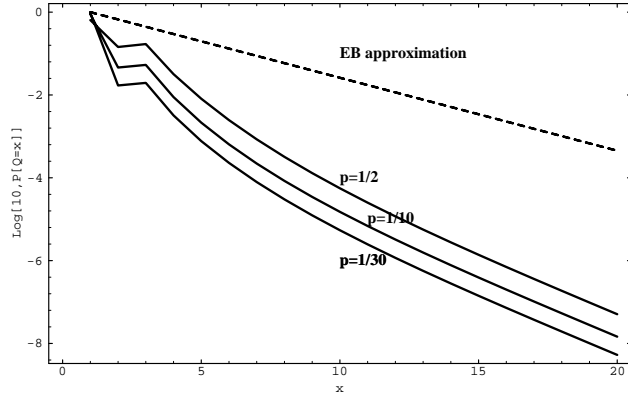This theorem is illustrated in the following numerical example.

FIGURE 5. Graph of $\log_{10} \mathbb{P}[Q = x]$ from Example 1 for different values of $T_{off}$ period parameter $p$.

**Example 2.2** *General on-off source.* Let (the modulating process) $B_t$ be a $\{0, 1\}$ valued process whose dynamics are described as follows. When in state zero (off), $B_t$ stays there for a geometrically distributed random time $T_{off}$, $\mathbb{P}[T_{off} = k] = (1-p)^{k-1} p$, $k \in \mathbb{N}$. When in state one (on) the process stays there for a generally distributed random time $T_{on}$ (independent of $T_{off}$). When in state zero the source is not producing anything $(Y(0) \equiv 0)$, and while in state one the source is producing i.i.d. arrivals with distribution $\mathbb{P}[Y(1) = 2] = 1 - \mathbb{P}[Y(1) = 0] = a$. Assume that the capacity of the server is $c = 1$. (Due to slightly simpler boundary conditions, all numerical examples in this paper were done for the recursion $Q_{t+1} = (Q_t - 1)^+ + A_t$; this recursion is asymptotically equivalent to (1.1).) From Theorem 2.1 it follows that as long as $T_{on}$ is long-tailed the equivalent bandwidth constant is independent of the distribution of $T_{on}$ and $T_{off}$ $(p)$. This is numerically illustrated for different values of $p$ in Figure 5; other parameters are taken to be $a = 2/5$, $\mathbb{P}[T_{on} \geq k] = k^{-3}, 1 \leq k \leq 80$, $\mathbb{P}[T_{on} \geq k] = 0, k > 80$. From the figure we can observe that the queue length probabilities are decreasing as $p$ decreases. This is intuitively obvious since off periods are getting larger. Also all the graphs eventually become parallel, as predicted by the previous theorem.

Although this example (as well as all the other examples in this paper) shows that the EB approximation is (too) conservative, it does not have to be so in general. In the context of the on-off source model we have observed that the EB approximation is too optimistic when the distribution of the on period decays faster than an exponential. Although we were not able to theoretically formulate this observation in greater generality, we believe it to be equally important. We illustrate this insight numerically. Take

$\mathbb{P}[T_{on} \geq k] = e^{-(k-1)^2}, k = 1, 2, 3, 4, 5, \mathbb{P}[T_{on} \geq k] = 0, k > 5$. Queue probabilities for $p = 1/10, a = 4/5$ are presented in Figure 6. We see that EB approximation is too optimistic. For more examples and details see [21]. Informally, if the input process "doesn't look exponential" the queue output is not exponential either.
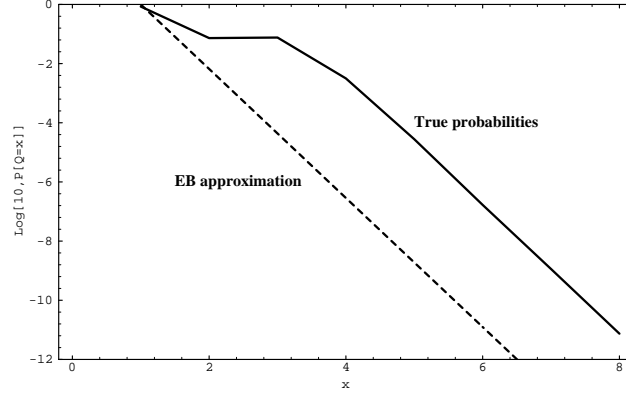


FIGURE 6. Comparison between the equivalent bandwidth approximation and the true probabilities for the case when the distribution of the on period decreases faster than an exponential.

## 3    Subexponential Arrivals

In this section our goal is to examine the asymptotics of the queue length distribution when the Cramér type conditions are replaced by subexponential assumptions. The two largest non Cramér families of distributions are long-tailed and subexponential distributions.

**Definition 3.1** *A distribution function F on $[0, \infty)$ is called* long-tailed *$(F \in \mathcal{L})$ if*

$$\lim_{x \to \infty} \frac{1 - F(x - y)}{1 - F(x)} = 1, \; y \in \mathbb{R}. \tag{1.8}$$

**Definition 3.2** *A distribution function F on $[0, \infty)$ is called* subexponential *$(F \in \mathcal{S})$ if*

$$\lim_{x \to \infty} \frac{1 - F^{*2}(x)}{1 - F(x)} = 2, \tag{1.9}$$

where $F^{*2}$ denotes the 2-nd convolution of $F$ with itself, i.e., $F^{*2}(x) = \int_{[0,\infty)} F(x-y)F(dy)$.

The class of subexponential distributions was first introduced by Chistakov [8]. The definition is motivated by the simplification of the asymptotic analysis of the convolution tails. Some examples of distribution functions in $\mathcal{S}$ are:

**(I)** the Pareto family

$$F(x) = 1 - (x - \beta + 1)^{-\alpha},$$

$x > \beta > 0, \alpha > 0$.

**(II)** the lognormal distribution

$$F(x) = \Phi\left(\frac{\log x - \mu}{\sigma}\right), \quad \mu \in \mathbb{R}, \sigma > 0,$$

where $\Phi$ is the standard normal distribution.

**(III)** Weibull distribution

$$F(x) = 1 - e^{-x^\beta},$$

for $0 < \beta < 1$.

**(IV)**

$$F(x) = e^{-x(\log x)^{-a}},$$

for $a > 0$. This class was proven to be subexponential in [33].

**(V)** Benktander Type I distribution [28]

$$F(x) = 1 - cx^{-a-1}x^{-b\log x}(a + 2b\log x),$$

$a > 0, b > 0$, and $c$ appropriately chosen.

**(V)** Benktander Type II distribution [28]

$$F(x) = 1 - cax^{-(1-b)}\exp\{-(a/b)x^b\},$$

$a > 0, 0 < b < 1$, and $c$ appropriately chosen.

The general relation between $\mathcal{S}$ and $\mathcal{L}$ is the following.

**Lemma 3.3** *(Athrey and Ney, [3])* $\mathcal{S} \subset \mathcal{L}$.

The following lemma [8] clearly shows that for long-tailed distributions Cramér type conditions are not satisfied.

**Lemma 3.4** *If $F \in \mathcal{L}$ then $(1 - F(x))e^{\alpha x} \to \infty$ as $x \to \infty$, for all $\alpha > 0$.*

An extensive treatment of subexponential distributions (and further references) can be found in Cline [11, 12].

Before we proceed any further, let us try to understand some of the basic properties of the sequence $\{X_n, n \geq 1\}$ of subexponentially distributed i.i.d. random variables. One of the main sample path characteristics of subexponential distributions follows from its definition [8], and that is

$$\mathbb{P}[X_1 + X_2 + \cdots + X_n > x] \sim n\mathbb{P}[X_1 > x], \qquad (1.10)$$

as $x \to \infty$. This means that a sum of subexponential random variables exceeds a large value $x$ by having one of them excede this value $x$; in terms of the appearance of the sample path of a sequence of subexponential random variables, we note that the sequence exhibits isolated peaks.
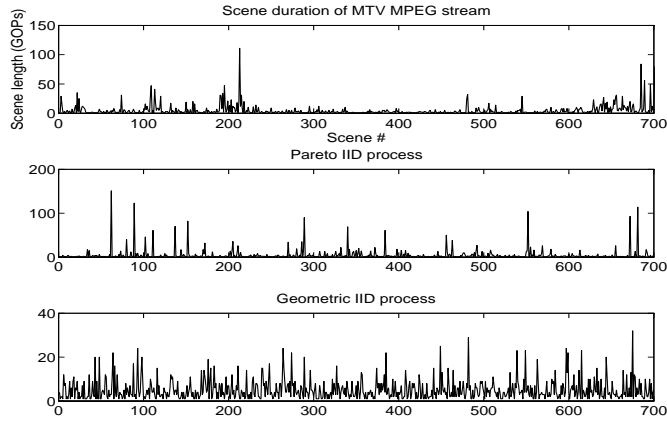


FIGURE 7. MPEG scene length duration (top); i.i.d. sample paths generated with the Pareto (middle) and geometric distribution (bottom).

Such a sample path behavior characterizes the scene lengths of video streams coded using the MPEG standard. Figure 7 shows a sequence of scene length durations (top), and for rough comparison, the sample paths generated by i.i.d. processes with Pareto (middle) and geometric distribution (bottom). Clearly, the scene length duration process has a subexponential character, as does the Pareto process, where the large peaks tend to be isolated in time, as suggested by (1.10). This is unlike the case of the geometrically distributed process. (For the description of MPEG data and the definition of scenes see [24].)

In terms of video traffic, subexponentiality can also manifest itself in the time-dependent (autocorrelation) structure. As shown in Figure 8, the autocorrelation function of MPEG video (17 streams multiplexed) matches the (subexponential) Pareto function $f(t) = \beta/t^\alpha$, for $\alpha = 0.513, \beta = 1.195$.
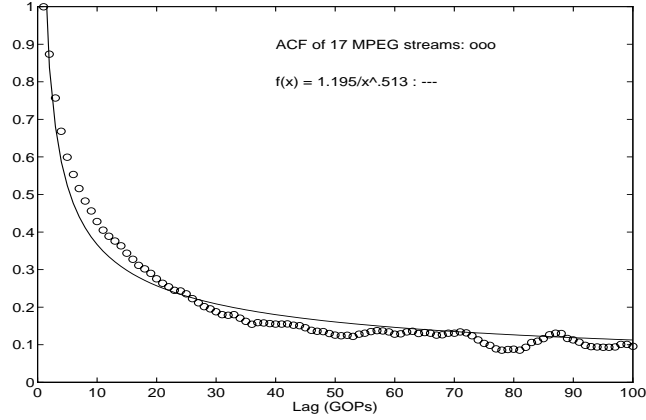
FIGURE 8. Modeling the autocorrelation function of MPEG video with an appropriate (subexponential) Pareto function.

In the next section, we will summarize some of the tools available for analyzing the queue behavior with subexponential arrivals.

### 3.1   Queueing Analysis

Assume that the queue increment process $X_t$ is a sequence of i.i.d. random variables with distribution function $F$, and $A_t$ is independent of $C_t$. Further, denote the *integrated tail* of $F$ as $\hat{F}(x) \stackrel{def}{=} \int_x^\infty [1 - F(t)]dt$, and define by $F_1(x) = m^{-1}(1 - \hat{F}(x))$, where $m = \hat{F}(0)$. Similarly, in the rest of the paper for any d.f. $G$, we define its corresponding $\hat{G}(x)$ and $G_1(x)$. Then the following result on the waiting time distribution asymptotics of the $GI/GI/1$ queue holds (see Veraverbeke [36]). Let $K$ be the d.f. of $A_t$.

**Theorem 3.5 (i)** $F_1 \in \mathcal{S} \Longleftrightarrow K_1 \in \mathcal{S}$ *and* $\lim_{x \to \infty} \frac{\hat{F}(x)}{\hat{K}(x)} = 1.$

**(ii)** *If* $K, K_1 \in \mathcal{S}$, *then*

$$\mathbb{P}[Q_t > x] \sim \frac{1}{\mathbb{E}C_t - \mathbb{E}A_t} \int_x^\infty \mathbb{P}[A_t > u]du, \text{ as } x \to \infty. \tag{1.11}$$

This theorem was first proved in [32]; in [36] the same result was shown using a random walk technique. Some of the first applications of long-tailed distributions in queueing theory were made by Cohen [13], and Borovkov [6] for functions of regular variations [26, 5]. Recent results on long-tailed and subexponential asymptotics of a $GI/GI/1$ are given in [1, 38]. (Also, in [1] further motivation is given for the application of long-tailed distributions to communication networks.)

The assumption that $K, K_1 \in \mathcal{S}$ in the theorem above can be replaced by an assumption on $K$ only. (Note that $K \in \mathcal{S}$ does not necessarily imply that $K_1 \in \mathcal{S}$.) This has been investigated in [28].

**Definition 3.6** $F \in \mathcal{S}^*$ *if*

$$\int_0^x \frac{\bar{F}(x-y)}{\bar{F}(x)} \bar{F}(y) dy \to 2m_F < \infty, \quad as \ \ x \to \infty,$$

*where* $m_F = \int_0^\infty yF(dy)$.

This class has the property that $\mathcal{S}^* \subset \mathcal{S}$, and that $F \in \mathcal{S}^* \Rightarrow F_1 \in \mathcal{S}$. Sufficient conditions for $F \in \mathcal{S}^*$ can be found in [29], where it was explicitly shown that lognormal, Pareto, and certain Weibull distributions are in $\mathcal{S}^*$.

An extension of Theorem 3.5 was investigated in [2]. In that paper the authors established the subexponential asymptotics of a Markov-modulated M/G/1 queue. However the constant of proportionality was left in a complex form. Full extension of Theorem 3.5 to Markov-modulated G/G/1 queues was given in [23] (preliminary results were reported in [22]), where it was proved that the queue length asymptotics are invariant under Markov modulation. A precise statement of this result follows.

Let $\{J_t\}$ be a stationary irreducible aperiodic Markov chain with a finite state space E (say with $N$ elements) and transition matrix $P$, and let $\{X_t\}$ be a sequence of real valued random variables. A stationary Markov process $\{(J_t, X_t)\}$ on $E \times \mathbb{R}$ whose transition distribution depends only on the first coordinate is called a Markov-modulated random walk (MMRW). This process is completely defined by its transition matrix measure $F_{ij}(B) = \mathbb{P}[J_1 = j, X_1 \in B | J_0 = i]$, and $F = \{F_{ij}\}$ (note that $\|F\| = F((-\infty, \infty)) = P$). Let $\{(J_t^r, X_t^r)\}$ denote the associated reversed process. This process is determined by the set of transition measures $F_{ij}^r(B) = \mathbb{P}[J_0 = j, X_1 \in B | J_1 = i]$, with $F^r = \{F_{ij}^r\}$ being the corresponding transition matrix measure.

Let $(J_t, A_t)$ and $(J_t, C_t)$ be two MMRWs such that $A_t$ and $C_t$ are conditionally independent given $J_{t-1}, J_t$; $\{A_t\}$ and $\{C_t\}$ are arrival and service processes, respectively. Let $K$ and $D$ be the corresponding transition measures for these MMRWs, i.e., $K = \{K_{ij}\} = \{\mathbb{P}[A_1 \in B, J_1 = j | J_0 = i]\}$, and $D = \{D_{ij}\} = \{\mathbb{P}[C_1 \in B, J_1 = j | J_0 = i]\}$; the reversed transition measure for the arrival process is $K^r = \{K_{ij}^r\} = \{\mathbb{P}[A_1 \in B, J_0 = j | J_1 = i]\}$, $B \in \mathcal{B}(\mathbb{R})$. For any (matrix) measure $H$, we denote $\bar{H}(x) = H(x, \infty)$. Then the following theorem holds [23].

**Theorem 3.7** *Let* $\lim_{x \to \infty} \overline{K^r}(x)/\bar{H}(x) = W$, *as* $x \to \infty$, $W = \{W_{ij}\}, W_{ij} \in [0, \infty)$, $H(x) \in \mathcal{L}, H_1(x) \in \mathcal{S}$ *(or* $H \in S^*$*), with at least one* $W_{ij} > 0$. *If* $\mathbb{E}C_t > \mathbb{E}A_t$, *and* $P$ *(=* $\|K\| = \|D\|$*) is irreducible and aperiodic, then,*

$$\frac{1}{\bar{H}(x)} \bar{Q}(x) \to \frac{1}{\mathbb{E}C_t - \mathbb{E}A_t} e\pi We, \quad as \ \ x \to \infty, \qquad (1.12)$$

where $\bar{Q}(x)$ *is a column vector with its* $i$*th component equal to* $\mathbb{P}[Q_t > x | J_t = i]$. *In particular,*

$$\mathbb{P}[Q_t > x] \sim \frac{1}{\mathbb{E}C_t - \mathbb{E}A_t} \int_x^\infty \mathbb{P}[A_t > u] du, \text{ as } x \to \infty. \quad (1.13)$$
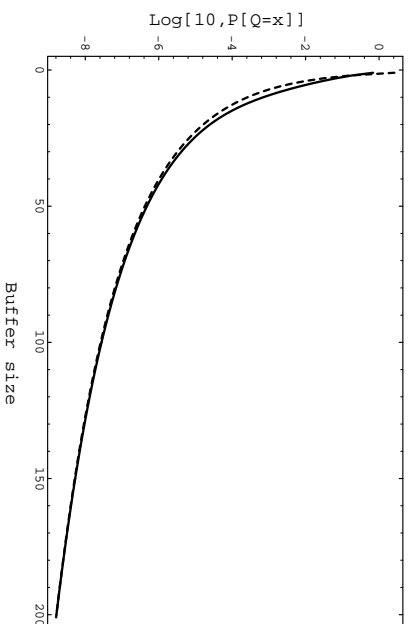


FIGURE 9. Graph of $\log_{10} \mathbb{P}[Q = i]$ versus buffer size $i$ from Example 2; the solid line represents the true probabilities, and the dashed line represents the approximation $2.603/i^{-4}$.

An illustration of the preceding theorem is given in the following numerical example.

**Example 3.8** Consider a constant server queue with $C_t = 1$ and two state (e.g. $\{0, 1\}$) Markov-modulated arrivals (source). The transition probabilities for the modulating Markov chain are $p_{01} = 1/3$, $p_{10} = 3/4$. When in state 0, the source is producing zero arrivals, and when in state 1, the source is producing (independently of the previous state) arrivals according to the distribution $\mathbb{P}[A_t = 0 | J_t = 1] = 0.327144$, $\mathbb{P}[A_t = 1 | J_t = 1] = 0$, and $\mathbb{P}[A_t = i | J_t = 1] = w/i^5$, $w = 18.220859$, $2 \le i \le 350$; $\rho_1 = \mathbb{E}[A_t | J_t = 1] = 3/2$. (Note that these are bounded arrivals.) Thus, according to the previous theorem, the queue length distribution is proportional to $1/i^4$, and the constant of proportionality is easily calculated to be $c = w\pi_1/(4(1 - \rho_1 \pi_1)) = w/7 = 2.603$. The comparison between the true probabilities and the approximation $c/i^4$ is shown in Figure 9.

*Stationary subexponentially correlated arrivals.* The models that we have seen in this section exhibit weak exponential autocorrelation structure and

dominant subexponential marginal distributions. For modeling subexponentially correlated arrivals in [23], we introduced the following class of processes. (These processes are a particular case of semi Markov processes [10].)

Consider a point process $T = \{T_0 \leq 0, T_n, n \geq 1\}$ such that $T_n - T_{n-1}, n \geq 1$ are i.i.d. with subexponential distribution function $F$. Further, let $J_n, n \geq 0$ be an irreducible aperiodic Markov Chain with finite state space $\{1, \ldots, K\}$, transition matrix $\{P_{ij}\}$, and stationary probability distribution $\pi_i, 1 \leq i \leq K$. In order to make this point process stationary (see [10], section 9.3), we choose the residual time at zero until the first jump to be distributed as an integrated tail of $F$, i.e., $F_1(t) = \mathbb{P}[T_1 \leq t] = m_F^{-1} \int_{0,t} \bar{F}(u)du, m_F = \mathbb{E}(T_n - T_{n-1})$.
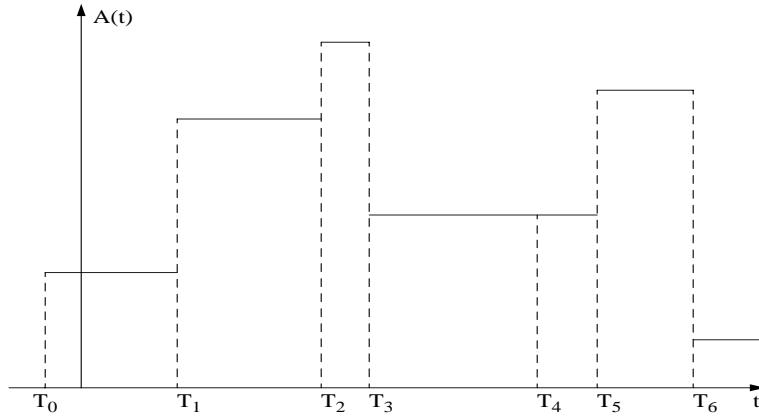


FIGURE 10. A possible realization of a Markov chain embedded into a renewal process.

Now we construct the following process:

$$A_t = J_n \quad \text{for} \quad T_n \leq t < T_{n+1}, \tag{1.14}$$

called a Markov Chain Embedded in a Stationary Subexponential Renewal Process (MCESSR). A typical sample path of this process is given in Figure 10. It is well known that under fairly general conditions, a Markov chain converges exponentially fast to its steady state distribution. These MCESSR processes have the characteristic that they, unlike finite state Markov chains, approach their steady state distributions with a subexponential rate. We illustrate this in the following example.

**Example 3.9** Let $F$ be a discrete distribution function with support $[1, 1000]$, $\mathbb{P}[T_2 - T_1 = 1] = 0.186532$, and $\mathbb{P}[T_2 - T_1 = i] = w/i^5, w = 22.028625$,

$2 \leq i \leq 1000$; choose a two state Markov chain with transition probabilities $p_{01} = 1/3$ and $p_{10} = 3/4$. Then, the functions $(d_{i,1}(t) \stackrel{def}{=} (\mathbb{P}_i[A_t = 1] - \pi_1)(\bar{F}_1(t)(\delta_{i_1} - \pi_1))^{-1}, i = 0, 1$, converge to one as $t \to \infty$, with subexponential rate. This can be clearly seen in Figure 11.
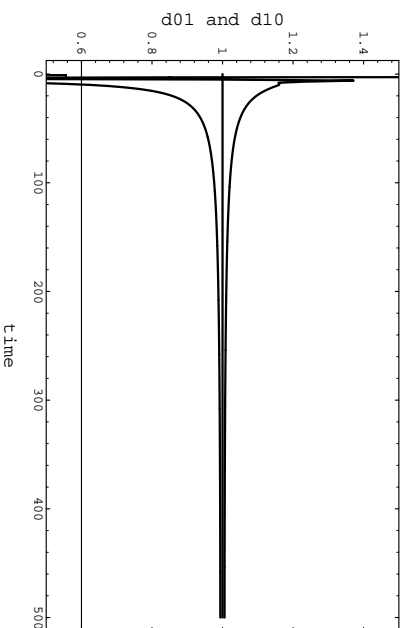


FIGURE 11. Functions $d_{i,1}(t) \stackrel{def}{=} (\mathbb{P}_i[A_t = 1] - \pi_1)(\bar{F}_1(t)(\delta_{i_1} - \pi_1))^{-1}, i = 0, 1$. The graph shows that $d_{i,1}(t) \to 1$ as $t \to 1$.

Another characteristic of these processes is that their autocorrelation functions, $R(t)$, are asymptotically proportional to the integrated tail of the sojourn time $T_n - T_{n-1}$ distribution, i.e. if $F, F_1 \in \mathcal{S}$, then

$$R(t) \sim \bar{F}_1(t),$$

as $t \to \infty$; this was formally proved in [23]. Combining these results, with Theorem 3.7, it was proven in the same paper that when the fluid flow queue is fed by these processes, its queue distribution is asymptotically proportional to its autocorrelation function, i.e.,

$$\mathbb{P}[Q > t] \sim^r R(t),$$

as $t \to \infty$. To the best of our knowledge, this was the first rigorous result relating the queue length distribution and the arrival process autocorrelation function.

## 4   Concluding Remarks

We have demonstrated that real-time traffic processes such as video traffic exhibit multiple time scale characteristics as well as subexponential first

and second order statistics. A network multiplexer that is loaded by these processes may manifest a distinct asymptotic behavior. We summarize recent results on evaluating the asymptotic behavior of a network multiplexer in the presence of subexponential and multiple time scale arrivals. It is left to identify in practice when some of the asymptotic techniques presented here can be applied to the design of efficient admission control policies in ATM based broadband networks.

*Acknowledgments:* The authors wish to thank the anonymous reviewer for his/her detailed list of editing suggestions.

## 5    References

[1] J. Abate, G.L. Choudhury, and W. Whitt. Waiting-time tail probabilities in queues with long-tail service-time distributions. *Queueing systems.*, 1994.

[2] S. Asmussen, L. F. Henriksen, and C. Klüppelberg. Large claims approximations for risk processes in a markovian environment. *Stochastic Processes and their Applications*, 54:29–43, 1994.

[3] K. B. Athreya and P. E. Ney. *Branching Processes*. Springer-Verlag, 1972.

[4] F. Baccelli and P. Bremaud. *Elements of Queueing Theory: Palm-Martingale Calculus and Stochastic Recurence*. Springer Verlag, 1994.

[5] N. H. Bingham, C. M. Goldie, and J. L. Teugels. *Regular Variation*. Cambridge University Press, Cambridge, 1987.

[6] A. A. Borovkov. *Stochastic Processes in Queueing Theory*. Springer-Verlag, 1976.

[7] C. S. Chang. Stability, queue length and delay of deterministic and stochastic queueing networks. *IEEE Transactions on Automatic Control*, 39:913–931, 1994.

[8] V. P. Chistakov. A theorem on sums of independent positive random variables and its application to branching random processes. *Theor. Probab. Appl.*, 9:640–648, 1964.

[9] G. L. Choudury, D. M. Lucantoni, and W. Whitt. Squeezing the most of ATM. To appear in *IEEE Trans. on Communications*, 1995.

[10] E. Cinlar. *Introduction to Stochastic Processes*. Prentice-Hall, 1975.

[11] D. B.H. Cline. Convolution tails, product tails and domains of attraction. *Probab. Th. Rel. Fields*, 72(1):529–557, 1986.

[12] D. B.H. Cline. Convolution of distributions with exponential and subexponential tails. *J. Austral. Math. Soc. (Series A)*, 43:347–365, 1987.

[13] J. W. Cohen. Some results on regular variation for distributions in queueing and fluctuation theory. *J. Appl. Probability*, 1973.

[14] C. Courcobetis and R. Weber. Buffer overflow asymptotics for a switch handling many traffic sources. Manuscript, December 1994.

[15] A. Elwalid, D. Heyman, T. V. Lakshman, D. Mitra, and A. Weiss. Fundamental bounds and approximations for atm multiplexers with applications to video teleconferencing. *IEEE Journal on Selected Areas in Communications*, 13(6):1004–1016, August 1995.

[16] A. I. Elwalid and D. Mitra. Effective bandwidth of general markovian traffic sources and admission control of high speed networks. *IEEE/ACM Trans. on Networking*, 1(3):329–343, June 1993.

[17] P. V. Glynn and W. Whitt. Logarithmic asymptotics for steady-state tail probabilities in a single-server queue. *Studies in Appl. Prob.*, 1994.

[18] R. Guerin, H. Ahmadi, and M. Nagshineh. Equivalent capacity and its application to bandwidth allocation in high-speed networks. *IEEE J. Select. Areas Commun.*, 9:968–981, 1991.

[19] P. R. Jelenković and A. A. Lazar. On the dependence of the queue tail distribution on multiple time scales of ATM multiplexers. In *Conference on Information Sciences and Systems*, pages 435–440, Baltimore, MD, March 1995. (www: http: //www.ctr.columbia.edu- /comet/publications).

[20] P. R. Jelenković and A. A. Lazar. Evaluating the queue length distribution of an ATM multiplexer with multiple time scale arrivals. In *Proceedings of INFOCOM'96*, San Francisco, California, March 1996, to appear.

[21] P. R. Jelenković and A. A. Lazar. Asymptotic properties of a discrete time queue with multiple time scale arrivals. *Submited to Queueing Systems*, 1995.

[22] P. R. Jelenković and A. A. Lazar. Subexponential asymptotics of a network multiplexer. In *Proceedings of the 33rd Annual Allerton Conference on Communication, Control, and Computing*, Urbana-Champaign, Illinois, October 1995.

[23] P. R. Jelenković and A. A. Lazar. Subexponential asymptotics of a markov-modulated G/G/1 queue. *Submitted to Journal of Appl. Prob.*, November 1995.

[24] P. R. Jelenković, A. A. Lazar, and N. Semret. Multiple time scales and subexponentiality in MPEG video streams. In *International IFIP-IEEE Conference on Broadband Communications*, April 1996, to appear.

[25] P. R. Jelenkovic and B. Melamed. Algorithmic modeling of TES processes. *IEEE Transactions on Automatic Control*, 40(7):1305–1312, July 1995.

[26] J. Karamata. Sur un mode de croissance réguilière des fonctions. *Mathematica (Cluj)*, 1930.

[27] F. P. Kelly. Effective bandwidths at multi-class queues. *Queueing Systems*, 9:5–16, 1991.

[28] C. Kluppelberg. Subexponential distributions and integrated tails. *J. Appl. Prob.*, 25:132–141, 1988.

[29] C. Kluppelberg. Subexponential distributions and characterizations of related classes. *Probability Theory and Related Fields*, 82:259, 1989.

[30] A. A. Lazar, G. Pacifici, and D. E. Pendarakis. Modeling video sources for real-time scheduling. *Multimedia Systems*, 1(6):253–266, 1994.

[31] R. M. Loynes. The stability of a queue with non-independent inter-arrival and service times. *Proc. Cambridge Philos Soc.*, 58:497–520, 1968.

[32] A.G. Pakes. On the tails of waiting-time distribution. *J. Appl. Probab.*, 12:555–564, 1975.

[33] E. J. G. Pitman. Subexponential distribution functions. *J. Austral. Math. Soc. Ser. A*, 29:337–347, 1980.

[34] D. Stoyan. *Comparison Methods for Queues and Other Stochastic Models*. John Wiley & Sons, 1983.

[35] D. Tse, R. Gallager, and J. Tsitsiklis. Statistical multiplexing of multiple time-scale markov stream. *IEEE, Selected Areas in Communications*, August 1995.

[36] N. Veraverbeke. Asymptotic behavior of wiener-hopf factors of a random walk. *Stochastic Proc. and Appl.*, 5:27–37, 1977.

[37] A. Weiss and A. Shwartz. *Large Deviations for Performance Analysis: Queues, Communications, and Computing*. New York: Chapman & Hall, 1995.

[38] E. Willekens and J. L. Teugels. Asymptotic expansion for waiting time probabilities in an M/G/1 queue with long-tailed service time. *Queueing Systems*, 10:295–312, 1992.