

ASYMPTOTIC RESULTS FOR MULTIPLEXING SUBEXPONENTIAL ON-OFF PROCESSES

PREDRAG R. JELENKOVIĆ,* ** AND
AUREL A. LAZAR,* *Columbia University*

Abstract

Consider an aggregate arrival process A^N obtained by multiplexing N on-off processes with exponential off periods of rate λ and subexponential on periods τ^{on} . As N goes to infinity, with $\lambda N \rightarrow \Lambda$, A^N approaches an $M/G/\infty$ type process. Both for finite and infinite N , we obtain the asymptotic characterization of the arrival process activity period.

Using these results we investigate a fluid queue with the limiting $M/G/\infty$ arrival process A_t^∞ and capacity c . When on periods are regularly varying (with non-integer exponent), we derive a *precise asymptotic behavior* of the queue length random variable Q_t^P observed at the beginning of the arrival process activity periods

$$\mathbb{P}[Q_t^P > x] \sim \Lambda \frac{r + \rho - c}{c - \rho} \int_{x/(r+\rho-c)}^{\infty} \mathbb{P}[\tau^{\text{on}} > u] du \quad x \rightarrow \infty,$$

where $\rho = \mathbb{E}A_t^\infty < c$; r ($c \leq r$) is the rate at which the fluid is arriving during an on period. The asymptotic (time average) queue distribution lower bound is obtained under more general assumptions regarding on periods than regular variation.

In addition, we analyse a queueing system in which one on-off process, whose on period belongs to a subclass of subexponential distributions, is multiplexed with independent exponential processes with aggregate expected rate $\mathbb{E}e_t$. This system is shown to be asymptotically equivalent to the same queueing system with the exponential arrival processes being replaced by their total mean value $\mathbb{E}e_t$.

Keywords: Non-Cramér type conditions; subexponential distributions; long-tailed distributions; long-range dependency; network multiplexer; fluid flow queue; $M/G/\infty$ queue

AMS 1991 Subject Classification: Primary 60J15
Secondary 60K25

1. Introduction

The problem of multiplexing on-off sources arises frequently as the basic model of contention in multimedia communication systems, as well as in some storage systems. More specifically, in modern multimedia communication networks, such as ATM, various calls are simultaneously established among the different source-destination pairs. These calls are usually discretized/packetized. An individual call/source can be either active, in which case it

Received 26 November 1996; revision received 6 July 1998.

* Postal address: Department of Electrical Engineering and Center for Telecommunications Research, Columbia University, New York, NY 10027, USA.

** Email address: predrag@ctr.columbia.edu

transmits a packet, or silent, i.e. there is no transmission. Therefore, each source can be viewed as an on-off process. Along their routes, sources that are simultaneously active share common network resources: bandwidth, buffer space, computational power, etc. The fundamental building blocks for sharing bandwidth and buffer space are network multiplexers (MUX). Sharing of common network resources may lead to quality of service (QOS) degradation for individual flows. Therefore, it is important to have computationally efficient algorithms for evaluating QOS measures under all possible traffic loads. A first step towards a satisfactory solution to this problem is understanding network MUX units in isolation. The most fundamental mathematical model of a network MUX is an infinite buffer queue loaded with multiplexed on-off arrival processes; the main QOS performance measure for this queueing system is the buffer occupancy probability distribution.

The problem of multiplexing dates back to [21, 50]. In [21], Cohen obtained a complete Laplace transform solution to this problem! More recently, he revisited this problem in [22]. However, inverting the Laplace transform is usually a very tedious process. Hence, investigating computationally tractable exact and approximate solution techniques are needed. For Markovian fluid on-off processes a thorough investigation of this problem was made in [3]. Many other results for multiplexing Markovian on-off processes followed. These led to the Equivalent Bandwidth theory for Markovian, or, in general, exponentially bounded arrival processes; extensive references can be found for example in [24, 25, 27].

Recently statistical analysis has increasingly shown that the traffic streams in modern broadband networks exhibit long-tailed/subexponential characteristics. For Ethernet traffic such results were examined in [41]. These statistical results have stimulated research in queueing analysis under the heavy-tailed (non-Cramér) assumptions. Queueing analysis with self-similar long-range dependent arrival processes appears in [24, 42, 44, 46, 48, 51]. Recently, long-tailed characteristics of the scene length distribution of MPEG video streams were explored in [29, 30, 35, 36].

Parallel to the modeling approach through self-similar long-range dependent processes, a more analytically tractable approach using fluid renewal type models in which renewal times are long-tailed has been explored in [2, 28]. Queueing results in these two papers rely on the classical result by Pakes [45] on the subexponential asymptotics of the waiting time distribution in a $GI/GI/1$ queue or on earlier work of Cohen [20] which considered a regularly varying $GI/GI/1$ queue.

The result of Pakes has been generalized to a Markov modulated setting [6, 34]. In [6] the subexponential asymptotics of a Markov modulated $M/G/1$ queue was investigated. Work in [34] further generalized these results to Markov modulated $G/G/1$ queues. In the same paper it was shown that a subexponential $GI/GI/1$ queue is invariant under Markov modulation. In other words, a subexponential Markov modulated $G/G/1$ queue has the same asymptotics as the corresponding $GI/GI/1$ queue. These results made possible the analysis of a subexponential semi-Markov fluid queue [34]. Further generalizations of the result in [34] to arrival processes with a more complex dependency structure were investigated in [7]. Asymptotic expansion refinements of Pakes' result can be found in [1, 53].

The analysis of a fluid queue in which more than one long-tailed process is multiplexed appears to be a very difficult problem. This is due to the fact that the renewal structure of an aggregate arrival process may be very complex, although the appearance of each individual process may be truly innocuous (like an on-off process). The complex autocorrelation structure of the aggregate process obtained by multiplexing long-tailed on-off processes has been examined in [28]. General bounds for multiplexing long-tailed fluid processes have been derived

in [15]. In [12] a limiting case of an infinite number of on-off processes with regularly varying on distribution has been investigated. In the same paper (see also [13]) a case of two processes, one of which had regularly varying on periods and the other had exponential on periods, has been solved. A similar scenario with intermediately regularly varying on periods has been examined in [49]. The literature does not explicitly give precise asymptotic results for the case of multiplexing two or more long-tailed processes.

From a mathematical point of view the new results in this paper were achieved through the combination of the theory of subexponential distributions, renewal theory, Karamata's theory and the utilization of sample path arguments. From an engineering standpoint this paper advances two important results. The first result intuitively states that when a process with subexponential characteristics, e.g. MPEG video, is multiplexed with a process that has exponential characteristics, e.g. voice, the contribution to the large buffer asymptotics of the exponential processes is reflected only through their mean values. This result suggests that, under appropriate conditions, for admission control of both VBR video and voice streams, the voice streams need to be characterized only by their mean values. The second result is an accurate approximation with low computational complexity of the large buffer probabilities of finitely many subexponential on-off processes. Besides accuracy, it is of special importance for engineering the MUX that this approximation has basically *negligible computational complexity*. To the best of our knowledge, this is the only result in literature of comparable computational complexity that is both proven theoretically and demonstrated experimentally as a good approximation for the buffer overflow probabilities with multiplexed long-tailed arrival streams. For the experimental verification of this asymptotic approximation see [31, 32, 33].

The rest of the paper is organized as follows. Section 2 contains necessary definitions and examples of long-tailed and subexponential distributions. In Section 3 we examine the aggregate arrival process obtained by multiplexing N independent and identically distributed on-off processes. For this process we derive the asymptotic relation between the distribution of its activity period and the distribution of on periods of individual processes. Using Karamata's theory for the case when on periods are regularly varying with non-integer exponent, we obtain a precise asymptotic behavior of the server overflow distribution during the arrival process activity period. Using these asymptotic relations, in Section 4, we derive several results for the fluid queue asymptotics of multiplexed long-tailed processes. The paper is concluded in Section 5.

2. Long-tailed and subexponential distributions

This section contains necessary definitions of long-tailed and subexponential distributions. For convenience, we give some basic results on these distributions in Appendix A.

Definition 2.1. A distribution function F on $[0, \infty)$ is called *long-tailed* ($F \in \mathcal{L}$) if

$$\lim_{x \rightarrow \infty} \frac{1 - F(x - y)}{1 - F(x)} = 1, \quad \forall y \in \mathbb{R}. \quad (2.1)$$

Definition 2.2. A distribution function F on $[0, \infty)$ is called *subexponential* ($F \in \mathcal{S}$) if

$$\lim_{x \rightarrow \infty} \frac{1 - F^{*2}(x)}{1 - F(x)} = 2, \quad (2.2)$$

where F^{*2} denotes the second 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 Chistyakov [14]. The definition is motivated by the simplification of the asymptotic analysis of convolution tails. One of the best known examples of distribution functions in \mathcal{S} (and \mathcal{L}) are functions of regular variation $\mathcal{R}_{-\alpha}$ (in particular Pareto family); $F \in \mathcal{R}_{-\alpha}$ if it is given by

$$F(x) = 1 - \frac{l(x)}{x^\alpha} \alpha \geq 0,$$

where $l(x) : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ is a function of slow variation, i.e. $\lim_{x \rightarrow \infty} l(\delta x)/l(x) = 1, \delta > 1$. These functions were invented by Karamata [37] (the main reference book is [11]). The other examples include lognormal and some Weibull distributions (see [39, 34]).

3. Analysis of the aggregate arrival process

This section consists of two parts. The first part is contained in Section 3.1, where we asymptotically relate the tail of the activity period distribution of an aggregate arrival process, obtained by multiplexing independent and identically distributed on-off processes, and the tail of the on period distribution of the individual on-off processes. The main results are given in Theorem 3.3 and Theorem 3.5. In the second part, Section 3.2, using Karamata’s theory we derive the asymptotic behavior of the distribution of the queue increment during the arrival process activity period. In Section 4, these results will be used to obtain asymptotic queueing results.

More formally, consider two independent sequences of i.i.d. random variables $\{\tau_n^{\text{off}}, n \geq 0\}, \{\tau_n^{\text{on}}, n \geq 0\}, \tau_0^{\text{off}} = \tau_0^{\text{on}} = 0$. Define a point process $T_n^{\text{off}} \stackrel{\text{def}}{=} \sum_{i=0}^n (\tau_i^{\text{off}} + \tau_i^{\text{on}}), n \geq 0$; this process will be interpreted as representing the beginnings of off periods in an on-off process. Further, define an on-off process a_t with rate r , as

$$a_t = r \text{ if } T_n^{\text{off}} - \tau_n^{\text{on}} \leq t < T_n^{\text{off}}, \quad n \geq 1,$$

and $a_t = 0$, otherwise. For the rest of the paper, unless otherwise specified, we will assume that τ_n^{off} is exponentially distributed with parameter λ , i.e., $\mathbb{P}[\tau_n^{\text{off}} > t] = e^{-\lambda t}, t \geq 0$. Also, τ_n^{on} is assumed to have a finite mean. Steady state probabilities of this process are given as $\pi_0 = \lim_{t \rightarrow \infty} \mathbb{P}[a_t = 0] = 1/(1 + \lambda \mathbb{E}\tau^{\text{on}}) = 1 - \pi_1$, where $\pi_1 = \lim_{t \rightarrow \infty} \mathbb{P}[a_t = r]$. Let $A^N = \sum_1^N a^i$, be an aggregate arrival process obtained by multiplexing N independent and identically distributed on-off processes $a^i, 1 \leq i \leq N$.

3.1. Asymptotic behavior of the aggregate process activity period

The central idea of this section is to relate the rate of convergence of $\mathbb{P}[a_t = 0]$ to its steady state and the tail of its on period (Section 3.1.1). Then, the rate of convergence for the aggregate process is easily computable from $\mathbb{P}[A_t^N = 0] = (\mathbb{P}[a_t = 0])^N$; from here one can refer back and relate $\mathbb{P}[A_t^N = 0]$ to the tail of its activity period distribution (Sections 3.1.2 and 3.1.3).

3.1.1. *Single on-off process: convergence to steady state.* Let us now investigate the speed of convergence of $\mathbb{P}[a_t = 0]$ to its steady state. For that reason we define a transient function

$F_{tr}(t) \stackrel{\text{def}}{=} \mathbb{P}[a_t = 1]/\pi_1$; notice that $F_{tr}(0) = 0$ and $\lim_{t \rightarrow \infty} F_{tr}(t) = 1$. Therefore, if $F_{tr}(t)$ is monotonic, it will represent a proper probability distribution function, in which case we call it the transient probability measure. To simplify the notation, throughout the rest of the paper, for any distribution F we define its tail function as $\bar{F}(x) \stackrel{\text{def}}{=} 1 - F(x)$; in addition, if F has a finite mean $m \stackrel{\text{def}}{=} \int_0^\infty \bar{F}(x) dx$, we denote its integrated tail distribution as $F_1(x) \stackrel{\text{def}}{=} m^{-1} \int_0^x \bar{F}(u) du$. In the following theorem we use a class of subexponential distributions named \mathcal{S}_d . For all practical purposes $F \in \mathcal{S}_d$ is virtually the same as $F \in \mathcal{S}$. A precise definition of \mathcal{S}_d is given in Definition A.1.

Theorem 3.1. *Let F be a distribution of τ^{on} with $\mathbb{E}\tau^{\text{on}} < \infty$ and let τ^{off} be exponentially distributed with parameter λ . Then, the Laplace–Stieltjes (LS) transform of the transient function is given as*

$$\tilde{F}_{tr}(s) = s \int_0^\infty e^{-st} F_{tr}(t) dt = \frac{(1 + \lambda \mathbb{E}\tau^{\text{on}})F_1(s)}{1 + \lambda \mathbb{E}\tau^{\text{on}} F_1(s)}. \tag{3.1}$$

If $F_1 \in \mathcal{S}$ and $\lambda \mathbb{E}\tau^{\text{on}} < 1$, then

$$\bar{F}_{tr}(t) \sim \frac{1}{1 + \lambda \mathbb{E}\tau^{\text{on}}} \bar{F}_1(t) \quad \text{as } t \rightarrow \infty. \tag{3.2}$$

If $F \in \mathcal{S}_d$ and $\lambda \mathbb{E}\tau^{\text{on}} < 1$, then (3.2) holds, and the density function $f_{tr}(t) \stackrel{\text{def}}{=} dF_{tr}(t)/dt$ satisfies

$$f_{tr}(t) \sim \frac{1}{1 + \lambda \mathbb{E}\tau^{\text{on}}} \frac{\bar{F}(t)}{\mathbb{E}\tau^{\text{on}}} \quad \text{as } t \rightarrow \infty. \tag{3.3}$$

Proof. Equation (3.1) follows directly from equation (2.1.3) in [21], and the observation that $\tilde{F}_1(s) = (1 - \tilde{F}(s))/(s\mathbb{E}\tau^{\text{on}})$. The asymptotic relation in (3.2) is a direct consequence of Theorem A.1, and the fact that for all $\lambda \mathbb{E}\tau^{\text{on}} < 1$ equation (3.1) implies

$$\bar{F}_{tr}(t) = (1 + \lambda \mathbb{E}\tau^{\text{on}}) \sum_{n=1}^\infty (-\lambda \mathbb{E}\tau^{\text{on}})^{n-1} \bar{F}_1^{*n}(t). \tag{3.4}$$

Similarly, equation (3.3) follows by differentiation of the equation above and Theorem A.2. This finishes the proof of the theorem.

For the remainder of this paper, it is of special interest to find sufficient conditions under which the transient function $F_{tr}(t)$ is a proper distribution function, i.e. monotonic. One set of sufficient conditions, for the case when off periods are large, is given in the following theorem. This is typically satisfied when there are a large number of processes with a small average arrival rate.

Theorem 3.2. *Let F be a distribution of τ^{on} with $\mathbb{E}\tau^{\text{on}} < \infty$ and let τ^{off} be exponentially distributed with parameter λ . For any fixed $F \in \mathcal{S}_d$, there exists a $\lambda_0 > 0$, such that for all $\lambda < \lambda_0$, $F_{tr}(t) \equiv F_{tr}^\lambda(t)$ is a probability distribution function, i.e., $f_{tr}(t) \geq 0, t \geq 0$.*

Proof. Differentiation of (3.4) gives

$$\begin{aligned} f_{\text{tr}}(t) &= (1 + \lambda \mathbb{E}\tau^{\text{on}}) \sum_{n=1}^{\infty} (-\lambda \mathbb{E}\tau^{\text{on}})^{n-1} f^{\otimes n}(t) \\ &= (1 + \lambda \mathbb{E}\tau^{\text{on}}) \left(f(t) - \lambda \mathbb{E}\tau^{\text{on}} \sum_{n=2}^{\infty} (-\lambda \mathbb{E}\tau^{\text{on}})^{n-2} f^{\otimes n}(t) \right) \\ &\geq (1 + \lambda \mathbb{E}\tau^{\text{on}}) \left(f(t) - \lambda \mathbb{E}\tau^{\text{on}} \sum_{n=2}^{\infty} (\lambda \mathbb{E}\tau^{\text{on}})^{n-2} f^{\otimes n}(t) \right), \end{aligned}$$

where $f(t) = \bar{F}(t)/\mathbb{E}\tau^{\text{on}}$. Now, by applying Lemma A.5(i), for $\epsilon > 0$,

$$\mathbb{E}\tau^{\text{on}} \sum_{n=2}^{\infty} (\lambda \mathbb{E}\tau^{\text{on}})^{n-2} f^{\otimes n}(t) \leq \frac{\mathbb{E}\tau^{\text{on}} C_{\epsilon} (1 + \epsilon)^2 f(t)}{1 - (1 + \epsilon)\lambda \mathbb{E}\tau^{\text{on}}} \stackrel{\text{def}}{=} \frac{C'_{\epsilon} f(t)}{1 - (1 + \epsilon)\lambda \mathbb{E}\tau^{\text{on}}},$$

for some $C_{\epsilon}, C'_{\epsilon} > 0$, and all $\lambda < 1/(\mathbb{E}\tau^{\text{on}}(1 + \epsilon))$. Therefore,

$$f_{\text{tr}}(t) \geq (1 + \lambda \mathbb{E}\tau^{\text{on}}) f(t) \left(1 - \frac{\lambda C'_{\epsilon}}{1 - (1 + \epsilon)\lambda \mathbb{E}\tau^{\text{on}}} \right) \geq 0,$$

for all $\lambda \leq 1/(\mathbb{E}\tau^{\text{on}}(1 + \epsilon) + C'_{\epsilon})$. This concludes the proof of the theorem.

3.1.2. Finite number of processes. In this section we consider an aggregate process A_t^N that is obtained by multiplexing N independent on-off processes, i.e. $A_t^N = \sum_{i=1}^N a_t^i$, where a_t^i are independent and identically distributed on-off processes. Note that the indicator process $\mathbf{1}(A_t^N = 0)$ is an on-off process with exponentially distributed off periods with parameter $N\lambda$. Let $\{I_n^{N,\text{off}}, I_n^{N,\text{on}}, n \geq 1\}$ be the lengths of the n th off and on periods in the indicator process $\mathbf{1}(A_t^N = 0)$, respectively. In the following theorem we characterize asymptotically the tail of the distribution function of $I_n^{N,\text{on}}$. Observe that the steady state probability of the aggregate process being in state 0 is given by $\Pi_0 = \pi_0^N$.

Theorem 3.3. Assume that $\lambda < \lambda_0$, where λ_0 is the same as in Theorem 3.2 (or that $\mathbb{P}[a_t = 0]$ is monotonic, and $\lambda \mathbb{E}\tau^{\text{on}} < 1$). If $F \in \mathcal{S}_d$, then

$$\int_t^{\infty} \mathbb{P}[I^{N,\text{on}} > u] du \sim (1 + \lambda \mathbb{E}\tau^{\text{on}})^{N-1} \int_t^{\infty} \mathbb{P}[\tau^{\text{on}} > u] du \quad \text{as } t \rightarrow \infty, \tag{3.5}$$

and

$$\mathbb{P}[I^{N,\text{on}} > t] \sim (1 + \lambda \mathbb{E}\tau^{\text{on}})^{N-1} \mathbb{P}[\tau^{\text{on}} > t] \quad \text{as } t \rightarrow \infty. \tag{3.6}$$

Proof. From the definition of $F_{\text{tr}}(t)$ and Theorem 3.1 it follows that

$$\mathbb{P}[a_t^i = 0] = \pi_0(1 + \pi_1 \bar{F}_1(t) + o(\bar{F}_1(t))), \quad 1 \leq i \leq N,$$

where $\pi_0 = 1 - \pi_1 = 1/(1 + \lambda \mathbb{E}\tau^{\text{on}})$, and F_1 is the integrated tail distribution of F . This implies that

$$\mathbb{P}[A_t = 0] = \mathbb{P}[a_t^i = 0]^N = \pi_0^N(1 + N\pi_1 \bar{F}_1(t) + o(\bar{F}_1(t))) \quad \text{as } t \rightarrow \infty.$$

Therefore, the transient function $F_{tr,N}(t)$ satisfies the following asymptotics

$$\bar{F}_{tr,N}(t) \stackrel{\text{def}}{=} 1 - \frac{\mathbb{P}[A_t^N = 1]}{1 - \pi_0^N} = \frac{\mathbb{P}[A_t^N = 0] - \pi_0^N}{1 - \pi_0^N} \sim \frac{\pi_0^N N \pi_1}{1 - \pi_0^N} \bar{F}_1(t) \quad \text{as } t \rightarrow \infty. \quad (3.7)$$

By Theorem 3.2, for all $\lambda < \lambda_0$, $F_{tr}(t)$ is a distribution function, implying that $\mathbb{P}[a_t^i = 0]$ and $\mathbb{P}[A_t^N = 0]$ are monotonic, which further implies that $F_{tr,N}(t)$ is a probability distribution function. Here, let $F_{N,1}$ be the integrated tail distribution of $I^{N,\text{on}}$ and $\tilde{F}_{N,1}$ be its LS transform. Then, as in (3.1), we obtain

$$\tilde{F}_{tr,N}(s) = \frac{(1 + N\lambda\mathbb{E}I^{N,\text{on}})\tilde{F}_{N,1}(s)}{1 + N\lambda\mathbb{E}I^{N,\text{on}}\tilde{F}_{N,1}(s)}.$$

After simple algebra, it follows that

$$\tilde{F}_{N,1}(s) = \frac{\pi_0^N \tilde{F}_{tr,N}(s)}{1 - (1 - \pi_0^N)\tilde{F}_{tr,N}(s)},$$

or, in the time domain

$$\bar{F}_{N,1}(t) = \pi_0^N \sum_{n=1}^{\infty} (1 - \pi_0^N)^{n-1} \overline{F_{tr,N}^{*n}}(t). \quad (3.8)$$

Now, $F \in \mathcal{S}_d$ implies $F_1 \in \mathcal{S}$ (Theorem 1.1, [40]), which by (3.7), and Lemma A.3(ii), yields $F_{tr,N} \in \mathcal{S}$. Hence, by Theorem A.1 and (3.8) it follows that

$$\bar{F}_{N,1}(t) \sim \pi_0^{-N} \bar{F}_{tr,N}(t) \sim \frac{N\pi_1}{1 - \pi_0^N} \bar{F}_1(t) \quad \text{as } t \rightarrow \infty, \quad (3.9)$$

where the second asymptotic relation follows from (3.7). By substituting $\mathbb{E}I^{N,\text{on}} = (1/\pi_0^N - 1)/(N\lambda)$ in (3.9) we obtain

$$\int_t^{\infty} \mathbb{P}[I^{N,\text{on}} > u] du \sim (1 + \lambda\mathbb{E}\tau^{\text{on}})^{N-1} \int_t^{\infty} \mathbb{P}[\tau^{\text{on}} > u] du \quad \text{as } t \rightarrow \infty,$$

which proves (3.5).

In order to prove (3.6) we first determine the asymptotic behavior of $f_{tr,N}$

$$\begin{aligned} f_{tr,N}(t) &\stackrel{\text{def}}{=} \frac{d}{dt} \frac{\mathbb{P}[A_t^N = 1]}{1 - \pi_0^N} \\ &= \frac{d}{dt} \frac{1 - \mathbb{P}[a_t = 0]^N}{1 - \pi_0^N} \\ &= \frac{d}{dt} \frac{1 - (1 - \pi_1 F_{tr}(t))^N}{1 - \pi_0^N} \\ &= \frac{N}{1 - \pi_0^N} (1 - \pi_1 F_{tr}(t))^{N-1} \pi_1 f_{tr}(t) \\ &\sim \frac{\pi_0^N N \pi_1}{1 - \pi_0^N} \frac{\bar{F}(t)}{\mathbb{E}\tau^{\text{on}}} \quad \text{as } t \rightarrow \infty, \end{aligned} \quad (3.10)$$

where the last asymptotic relation follows from (3.3). Next, by taking a derivative in (3.8) we obtain

$$\bar{F}_N(t) = \mathbb{E}I^{N,\text{on}} \pi_0^N \sum_{n=1}^{\infty} (1 - \pi_0^N)^{n-1} f_{\text{tr},N}^{\otimes n}(t), \tag{3.11}$$

where $F_N(t)$ is the distribution function of $I^{N,\text{on}}$. Finally, by applying Theorem A.2 we get

$$\bar{F}_N(t) \sim \mathbb{E}I^{N,\text{on}} \pi_0^{-N} f_{\text{tr},N}(t) \text{ as } t \rightarrow \infty,$$

which together with (3.10) yields (3.6).

3.1.3. Infinite number of on-off processes. In this subsection we analyse the limiting case of an infinite number of on-off sources. First, we show that the aggregate process A_t^N converges in distribution to an $M/G/\infty$ process A_t^∞ which we define as follows. Let $T_n, n \geq 0, T_0 = 0$, be a Poisson process with rate Λ . Define $A_t^\infty = \sum_{n=1}^{\infty} r \mathbf{1}(T_n \leq t < T_n + \tau_n^{\text{on}}), r > 0$. Then the following theorem holds.

Theorem 3.4. *If $\mathbb{E}\tau_n^{\text{on}} < \infty$, and $\lambda N \rightarrow \Lambda$ as $N \rightarrow \infty$, then*

$$A_t^N \xrightarrow{d} A_t^\infty \text{ as } N \rightarrow \infty, \tag{3.12}$$

where \xrightarrow{d} symbolizes convergence in distribution.

Proof. It is enough to prove that the beginnings of on periods in the process A_t^N converge to a Poisson process with rate Λ . This follows from a classical result on multiplexing a large number of renewal processes [16, 23].

Lemma 3.1. *The transient probability of the arrival process A_t^∞ being silent is given by*

$$\mathbb{P}[A_t^\infty = 0] = \exp \left\{ -\Lambda \int_0^t \mathbb{P}[\tau^{\text{on}} > u] du \right\}. \tag{3.13}$$

Furthermore, if $\mathbb{E}\tau^{\text{on}} < \infty$, then $\lim_{t \rightarrow \infty} \mathbb{P}[A_t^\infty = 0] = \exp\{-\Lambda \mathbb{E}\tau^{\text{on}}\}$.

Proof. Follows from Theorem 2.2 in [21].

Remark. Observe that A_t^∞ represents the number of customers in service in an $M/G/\infty$ queue in which the customer service requirement has the same distribution as τ^{on} and the arrival rate is Λ . (For recent asymptotic results on $M/G/\infty$ processes see [47].)

Note that Theorem 3.3 implies that

$$\lim_{\substack{N \rightarrow \infty \\ \lambda N \rightarrow \Lambda}} \lim_{t \rightarrow \infty} \frac{\mathbb{P}[I^{N,\text{on}} > t]}{\mathbb{P}[\tau^{\text{on}} > t]} = \exp\{\Lambda \tau^{\text{on}}\}. \tag{3.14}$$

However, this does not necessarily imply that we can interchange the limit and derive the asymptotics of the activity period $I^{\infty,\text{on}}$ in the A_t^∞ process. The following result gives the asymptotic characterization of $I^{\infty,\text{on}}$ and indeed shows that the limits in (3.14) can be interchanged.

Theorem 3.5. *The asymptotics of the distribution of $I^{\infty,\text{on}}$ and its integrated tail are related as follows:*

(i) *If $F_1 \in \mathcal{S}$, then*

$$\int_t^\infty \mathbb{P}[I^{\infty,\text{on}} > u] du \sim \exp\{\Lambda \mathbb{E}\tau^{\text{on}}\} \int_t^\infty \mathbb{P}[\tau^{\text{on}} > u] du \quad \text{as } t \rightarrow \infty.$$

(ii) *If in addition $F \in \mathcal{S}_d$, then*

$$\mathbb{P}[I^{\infty,\text{on}} > t] \sim \exp\{\Lambda \mathbb{E}\tau^{\text{on}}\} \mathbb{P}[\tau^{\text{on}} > t] \quad \text{as } t \rightarrow \infty.$$

Remark. For the case of τ^{on} being regularly varying, $\mathbb{P}[\tau^{\text{on}} > t] = l(t)/t^\alpha$, $1 < \alpha < 2$, this result was obtained in [12] where Karamata’s Tauberian/Abelian theorems was used to asymptotically relate $I^{\infty,\text{on}}$ and τ^{on} .

Proof. As before, we define a transient probability distribution function

$$F_{\text{tr},\infty}(t) \stackrel{\text{def}}{=} \frac{\mathbb{P}[A_t^\infty > 0]}{1 - \exp\{-\Lambda \mathbb{E}\tau^{\text{on}}\}}.$$

Observe that this is a proper probability distribution function, i.e. it is monotonically increasing from $F_{\text{tr},\infty}(0) = 0$ to $\lim_{t \rightarrow \infty} F_{\text{tr},\infty}(t) = 1$. Next, by using Lemma 3.1 we derive

$$\begin{aligned} \bar{F}_{\text{tr},\infty}(t) &= 1 - \frac{1 - \exp\left\{-\Lambda \int_0^t \mathbb{P}[\tau^{\text{on}} > u] du\right\}}{1 - \exp\{-\Lambda \mathbb{E}\tau^{\text{on}}\}} \\ &= \frac{\exp\left\{-\Lambda \int_0^t \mathbb{P}[\tau^{\text{on}} > u] du\right\} \left(1 - \exp\left\{-\Lambda \int_t^\infty \mathbb{P}[\tau^{\text{on}} > u] du\right\}\right)}{1 - \exp\{-\Lambda \mathbb{E}\tau^{\text{on}}\}} \\ &\sim \frac{\Lambda \exp\{-\Lambda \mathbb{E}\tau^{\text{on}}\}}{1 - \exp\{-\Lambda \mathbb{E}\tau^{\text{on}}\}} \int_t^\infty \mathbb{P}[\tau^{\text{on}} > u] du \end{aligned} \tag{3.15}$$

as $t \rightarrow \infty$. Now, the Laplace transform of $F_{\text{tr},\infty}$ for any $s \in \mathbb{R}^+$ is given by

$$\begin{aligned} \tilde{F}_{\text{tr},\infty}(s) &= \int_0^\infty e^{-st} dF_{\text{tr},\infty}(t) \\ &= s \int_0^\infty e^{-st} F_{\text{tr},\infty}(t) dt \\ &= \frac{s}{1 - \exp\{-\Lambda \mathbb{E}\tau^{\text{on}}\}} \mathbb{E} \int_0^\infty e^{-st} \mathbf{1}(A_t^\infty > 0) dt \\ &= \frac{s}{1 - \exp\{-\Lambda \mathbb{E}\tau^{\text{on}}\}} \mathbb{E} \sum_{n=0}^\infty \int_{I_{n+1}^{\infty,\text{off}} + \sum_{i=0}^n (I_i^{\infty,\text{on}} + I_i^{\infty,\text{off}})}^{\sum_{i=0}^{n+1} (I_i^{\infty,\text{on}} + I_i^{\infty,\text{off}})} e^{-st} dt \\ &= \frac{s}{1 - \exp\{-\Lambda \mathbb{E}\tau^{\text{on}}\}} \mathbb{E} \sum_{n=0}^\infty \exp\left\{-s \left(I_{n+1}^{\infty,\text{off}} + \sum_{i=0}^n (I_i^{\infty,\text{on}} + I_i^{\infty,\text{off}})\right)\right\} \\ &\quad \times \frac{(1 - \exp\{-s I_{n+1}^{\infty,\text{on}}\})}{s} \\ &= \frac{\Lambda}{1 - \exp\{-\Lambda \mathbb{E}\tau^{\text{on}}\}} \frac{(1 - \mathbb{E} \exp\{-s I^{\infty,\text{on}}\})}{s + \Lambda - \Lambda \mathbb{E} \exp\{-s I^{\infty,\text{on}}\}}. \end{aligned}$$

By solving the last equation in $\mathbb{E} \exp\{-sI^{\infty, \text{on}}\}$ and putting $\gamma = 1 - \exp\{-\Lambda \mathbb{E} \tau^{\text{on}}\} (< 1)$, we obtain

$$\mathbb{E} \exp\{-sI^{\infty, \text{on}}\} = 1 - \frac{s\Lambda^{-1}\gamma \tilde{F}_{\text{tr}, \infty}(s)}{1 - \gamma \tilde{F}_{\text{tr}, \infty}(s)},$$

or equivalently

$$\frac{1 - \mathbb{E} \exp\{-sI^{\infty, \text{on}}\}}{s} = \frac{\Lambda^{-1}\gamma \tilde{F}_{\text{tr}, \infty}(s)}{1 - \gamma \tilde{F}_{\text{tr}, \infty}(s)}.$$

Observing that $(1 - \mathbb{E} \exp\{-sI^{\infty, \text{on}}\})/s$ is the LS transform of $\int_0^t \mathbb{P}[I^{\infty, \text{on}} > u] du$ we arrive at

$$\int_0^t \mathbb{P}[I^{\infty, \text{on}} > u] du = \Lambda^{-1} \sum_{n=1}^{\infty} \gamma^n (F_{\text{tr}, \infty})^{*n}(t), \tag{3.16}$$

which in combination with (3.15) and the same arguments as in the proof of Theorem 3.3 yields the conclusion (i) of this theorem.

The proof of (ii) can be obtained in a similar manner by first deriving the asymptotic behavior of $f_{\text{tr}}(t)$ and combining it with the derivative of (3.16); we omit the details.

3.2. Total server overflow during the activity period

Let $B_n, n \geq 1$, be a sequence of random variables representing the total amount of fluid that is brought to the system during the n th activity period, i.e., $B_n = \int_{t_n^b}^{t_n^e} A_t^{\infty} dt$, where t_n^b, t_n^e represent the beginning and end of the n th activity period, respectively. Further, define $D_{c,n} \stackrel{\text{def}}{=} B_n - cI_n^{\text{on}}, 0 < c \leq r$; note that $D_n \equiv D_{c,n}$ is a non-negative random variable. If we imagine that A_t^{∞} represents the rate at which the fluid is arriving to a fluid queue, and that c is the constant rate at which the queue drains, then D_n represents the queue increment during the n th activity period. In order to derive the queueing asymptotics, we first have to understand the asymptotic behavior of D_n . Unfortunately, this is a much more difficult task than the investigation of the asymptotic behavior of the activity period that we have done so far. For that reason we are forced to work under much more restrictive assumptions with distribution functions of regular variation. The method of proof for the following result will be through Karamata's Tauberian/Abelian theorems.

Theorem 3.6. *Consider an $M/G/\infty$ arrival process with on periods being regularly varying $\mathbb{P}[\tau^{\text{on}} > x] = l(x)/x^{\alpha}, \alpha > 1$, where α is non-integer. If $0 < c \leq r$, then*

$$\mathbb{P}[D_{c,n} > x] \sim \exp\{\Lambda \mathbb{E} \tau^{\text{on}}\} \mathbb{P}\left[\tau^{\text{on}} > \frac{x}{r + r\Lambda \mathbb{E} \tau^{\text{on}} - c}\right] \text{ as } x \rightarrow \infty. \tag{3.17}$$

Proof. Given in Appendix B.

Next, consider a stationary version of the arrival process

$$A_t^{\infty, s} = \sum_{-\infty < n < \infty} r \mathbf{1}(T_n \leq t < T_n + \tau_n^{\text{on}}),$$

where T_n is a stationary Poisson process with rate Λ . Given that at time $t = 0$, the arrival process is active ($A_t > 0$), denote by $D_{c(0)}$ the total queue increment since the beginning of the last activity period until time zero, i.e., $D_{c(0)} = \int_{t_0^b}^0 (A_t^\infty - c) dt$, $0 < c \leq r$, where t_0^b represents the beginning of the activity period that is still active at $t = 0$.

Now, by Theorem 4.3 of [4, p. 64], it follows that the process $\{T_n + \tau_n^{\text{on}}, -\infty < n < \infty\}$ is also a stationary Poisson process with the same rate Λ . Therefore, the process $A_t^{\infty,s}$ is reversible. This implies that $D_{c(0)}$ is equal in distribution to $\int_0^{t_0^e} (A_t^\infty - c) dt$, where t_0^e represents the end of the activity period that is active at $t = 0$. For simplicity reasons, we will refer to both of these variables by $D_{c(0)}$.

Conjecture 3.1. *Consider an $M/G/\infty$ arrival process with on periods being regularly varying, $\mathbb{P}[\tau^{\text{on}} > x] = l(x)/x^\alpha$, $\alpha > 1$, where α is non-integer. If $0 < c \leq r$, then*

$$\mathbb{P}[D_{c(0)} > x] \sim \frac{\Lambda \exp\{\Lambda \mathbb{E}\tau^{\text{on}}\}}{\exp\{\Lambda \mathbb{E}\tau^{\text{on}}\} - 1} \int_{x/(r+r\Lambda \mathbb{E}\tau^{\text{on}}-c)}^\infty \mathbb{P}[\tau^{\text{on}} > u] du \quad \text{as } x \rightarrow \infty. \quad (3.18)$$

Heuristics. We believe that the proof of this theorem can be done in the same spirit as the proof of Theorem 3.6. Unfortunately, this seems to be very technical and for that reason we do not attempt to provide a rigorous proof. Instead we give the following heuristics. From Theorem 3.6 it follows that

$$\mathbb{P}[D_c > x] \sim \mathbb{P}[I^{\text{on}}(r + r\Lambda \mathbb{E}\tau^{\text{on}} - c) > x] \quad \text{as } x \rightarrow \infty. \quad (3.19)$$

Based on this, one can expect that

$$\mathbb{P}[D_{c(0)} > x] \sim \mathbb{P}[I_{(0)}^{\text{on}}(r + r\Lambda \mathbb{E}\tau^{\text{on}} - c) > x] \quad \text{as } x \rightarrow \infty,$$

where $I_{(0)}^{\text{on}}$ is the residual activity time at time zero, which satisfies

$$\begin{aligned} \mathbb{P}[I_{(0)}^{\text{on}} > x] &= 1/\mathbb{E}I^{\text{on}} \int_x^\infty \mathbb{P}[I^{\text{on}} > u] du \\ &\sim 1/\mathbb{E}I^{\text{on}} e^{\Lambda \mathbb{E}\tau} \int_x^\infty \mathbb{P}[\tau^{\text{on}} > u] du \quad \text{as } x \rightarrow \infty, \end{aligned} \quad (3.20)$$

where (3.20) follows from Theorem 3.5. Finally, (3.18) follows by combining (3.20), (3.19), and $\mathbb{E}I^{\text{on}} = (\exp\{\Lambda \mathbb{E}\tau^{\text{on}}\} - 1)/\Lambda$.

4. Queuing analysis

We begin this section with a classical result on subexponential asymptotics of a $GI/GI/1$ queue. The result was obtained by Pakes 1975 (see also Veraverbeke for the random walk approach to this problem). For extensions of this result to Markov-modulated $M/G/1$ queues see [6], and to Markov-modulated $G/G/1$ queues (equivalently random walks) see [34]. Further extension of these results to more general arrival processes was obtained in [7].

Let $\{A_n, n \geq 0\}$ and $\{C_n, n \geq 0\}$ be two independent sequences of non-negative i.i.d. random variables that are driving a queueing process (Lindley’s recursion)

$$Q_{n+1} = (Q_n + A_n - C_n)^+, \quad n \geq 0, \quad (4.1)$$

where $q^+ = \max(0, q)$. According to the classical result of Loynes [43] under the stability condition $\mathbb{E}A_n < \mathbb{E}C_n$ this recursion admits a unique stationary solution, and for all initial conditions $\mathbb{P}[Q_n \leq x]$ converges to the stationary distribution $\mathbb{P}[Q \leq x]$. For the rest of this paper we will assume that all queueing systems under consideration are in their stationary regimes. Let G and G_1 represent the distribution and its integrated tail distribution for A_n , respectively.

Theorem 4.1 (Pakes) *If $G_1 \in \mathcal{S}$ (or $G \in \mathcal{S}_d$) and $\mathbb{E}A_n < \mathbb{E}C_n$, then*

$$\mathbb{P}[Q_n > t] \sim \frac{1}{\mathbb{E}C_n - \mathbb{E}A_n} \int_x^\infty \mathbb{P}[A_n > u] du \quad \text{as } t \rightarrow \infty.$$

4.1. Fluid queue: preliminaries

The physical interpretation for a fluid queue is that at any moment of time t , fluid is arriving to the system with rate a_t and is leaving the system with rate c_t . We call a_t and c_t the arrival and the service process, respectively. Then, the amount of fluid Q_t (also called queue length) evolves according to

$$dQ_t = (a_t - c_t) dt \quad \text{if } Q_t > 0, \text{ or } a_t > c_t, \tag{4.2}$$

and $dQ_t = 0$, otherwise. It is not very difficult to see that, starting from $Q_0 = 0$, the solution $Q_t, t \geq 0$, to (4.2) is given by

$$Q_t = \sup_{0 \leq u \leq t} \int_u^t (a_v - c_v) dv. \tag{4.3}$$

And if a_t and c_t are stationary, Q_t is equal in distribution to

$$\mathbb{P}[Q_t \leq x] = \mathbb{P} \left[\sup_{0 \leq u \leq t} W_u \leq x \right],$$

where $W_t \stackrel{\text{def}}{=} \int_{-t}^0 (a_u - c_u) du, t \geq 0$. Now, whenever the stability condition $\mathbb{E}a_t < \mathbb{E}c_t$ is satisfied (by Birkhoff's Strong Law of Large Numbers), $\mathbb{P}[Q_t \leq x]$ converges to a proper probability distribution, i.e.

$$\mathbb{P}[Q \leq x] \stackrel{\text{def}}{=} \lim_{t \rightarrow \infty} \mathbb{P}[Q_t \leq x] = \mathbb{P} \left[\sup_{0 \leq u < \infty} W_u \leq x \right].$$

Further, when the difference process $x_t \stackrel{\text{def}}{=} a_t - c_t$ is driven by a stationary and ergodic point process $\{T_n, -\infty < n < \infty\}$, i.e.

$$x_t = x_{T_n}, \quad t \in [T_n, T_{n+1}),$$

then the fluid queue process evolves as

$$Q_t = (Q_{T_n-} + (t - T_n)x_{T_n})^+, \quad t \in [T_n, T_{n+1}), \tag{4.4}$$

where $q^+ = \max(q, 0)$. From the recursion above, it is clear that the process Q_t is essentially the same as the $G/G/1$ workload process. Hence, by the fundamental stability theorem of Loynes (see Chapter 2 in [9]) there exists a unique stationary process $\{Q_t^s, -\infty < t < \infty\}$ ($\mathbb{P}[Q_t^s \leq x] = \mathbb{P}[Q \leq x]$) that satisfies (4.4) (or equivalently (4.2)). In the rest of the paper, whenever we refer to Q_t , we will actually mean Q_t^s . The existence and uniqueness of this stationary solution will be important for establishing the relation between the Palm queue probabilities and the time average probabilities $\mathbb{P}[Q_t \leq x]$.

4.2. Fluid queue with a single on-off process

Consider a fluid queue with capacity c and an on-off arrival process with on arrival rate r . In this subsection we assume that off periods are also general (not necessarily exponential). A general storage model in a two state random environment was investigated in [38]. Then, if we observe the queue at the beginning of on periods, the queue length Q_n^P evolves as follows (P stands for Palm probability [9]).

$$Q_{n+1}^P = (Q_n^P + (r - c)\tau_n^{\text{on}} - c\tau_n^{\text{off}})^+, \quad n \geq 0. \quad (4.5)$$

Recall that F and F_1 denote the distribution and the integrated tail distribution of τ^{on} .

Theorem 4.2. *If $r > c$, $(r - c)\mathbb{E}\tau_{\text{on}} < c\mathbb{E}\tau_{\text{off}}$, and $F_1 \in \mathcal{S}$ (or $F \in \mathcal{S}_d$), then*

$$\mathbb{P}[Q_n^P > x] \sim \frac{r - c}{c\mathbb{E}\tau_{\text{off}} - (r - c)\mathbb{E}\tau_{\text{on}}} \int_{x/(r-c)}^{\infty} \mathbb{P}[\tau^{\text{on}} > u] du \quad \text{as } x \rightarrow \infty. \quad (4.6)$$

Proof. By defining $A_n = (r - c)\tau_n^{\text{on}}$ and $C_n = c\tau_n^{\text{off}}$ the theorem follows immediately from Theorem 4.1.

4.2.1. Time averages. Here, we will compute queue time averages based on the queue Palm probabilities computed in Theorem 4.2. For this we need a stationary version a_t^s of the on-off arrival process a_t . Let $T_n^{\text{on}}, -\infty < n < \infty$, be a stationary point process that represents the beginnings of the on-off periods, with a convention that $T_0^{\text{on}} < 0 \leq T_1^{\text{on}}$. Then, according to [48], the random variable T_0^{on} can be represented as $-T_0^{\text{on}} = B(\tau_{(0)}^{\text{off}} + \tau_{(0)}^{\text{on}}) + (1 - B)\tau_{(0)}^{\text{on}}$, where the random variables $B, \tau_{(0)}^{\text{on}}, \tau_{(0)}^{\text{off}}$ are independent of $\{\tau_n^{\text{on}}, \tau_n^{\text{off}}, n \leq -1\}$, $\tau_{(0)}^{\text{off}}, B$ is a Bernoulli random variable with $\mathbb{P}[B = 0] = 1 - \mathbb{P}[B = 1] = \mathbb{E}\tau^{\text{on}} / (\mathbb{E}\tau^{\text{on}} + \mathbb{E}\tau^{\text{off}})$, and $\tau_{(0)}^{\text{on}}, \tau_{(0)}^{\text{off}}$ are distributed as integrated tail distributions of $\tau^{\text{on}}, \tau^{\text{off}}$, respectively. Furthermore, the net increment of the load that comes to the queue in the interval $[T_0, 0]$ is given by the following equation

$$\int_{T_0^{\text{on}}}^0 (a_t^s - c) dt = B[(r - c)\tau_{(0)}^{\text{on}} - c\tau_{(0)}^{\text{off}}] + (1 - B)(r - c)\tau_{(0)}^{\text{on}} \quad (4.7)$$

Before we present our result let us state the following well known lemma on long-tailed distributions.

Lemma 4.1. *Let X and Y be two independent non-negative random variables. If $X \in \mathcal{L}$, then*

$$\mathbb{P}[X - Y > t] \sim \mathbb{P}[X > t] \quad \text{as } t \rightarrow \infty.$$

Proof. Follows easily from the definition of \mathcal{L} .

Theorem 4.3. *If $r > c$, $(r - c)\mathbb{E}\tau_{\text{on}} < c\mathbb{E}\tau_{\text{off}}$, and $F_1 \in \mathcal{S}$ (or $F \in \mathcal{S}_d$), then*

$$\mathbb{P}[Q_t > x] \sim \mathbb{P}[Q_n^P > x] + \frac{1}{\mathbb{E}\tau^{\text{off}} + \mathbb{E}\tau^{\text{on}}} \int_{x/(r-c)}^{\infty} \mathbb{P}[\tau^{\text{on}} > u] du \quad (4.8)$$

$$\sim K \int_{x/(r-c)}^{\infty} \mathbb{P}[\tau^{\text{on}} > u] du \quad \text{as } x \rightarrow \infty, \quad (4.9)$$

where

$$K = \frac{r - c}{c\mathbb{E}\tau_{\text{off}} - (r - c)\mathbb{E}\tau_{\text{on}}} + \frac{1}{\mathbb{E}\tau_{\text{off}} + \mathbb{E}\tau_{\text{on}}}. \tag{4.10}$$

Remarks. (i) This theorem improves on known results in [15, 48] that were obtained under the assumption of τ^{on} being regularly varying.

(ii) The following proof can be carried out to establish the relationship between the Palm and time averages in much more general settings like semi-Markov fluid queues.

Proof. Let $\{Q_t, -\infty < t < \infty\}$ be a unique stationary solution to (4.5). Then, by using equation (4.7), and the independence of B of $Q_{T_0}, \tau_{(0)}^{\text{off}}, \tau_{(0)}^{\text{on}}, \tau_{(0)}^{\text{off}}$, we obtain

$$\begin{aligned} \mathbb{P}[Q_0 > x] &= \mathbb{P}[Q_0 > x, B = 1] + \mathbb{P}[Q_0 > x, B = 0] \\ &= \mathbb{P}[Q_{T_0} + \tau_0^{\text{on}}(r - c) - c\tau_{(0)}^{\text{off}} > x, B = 1] + \mathbb{P}[Q_{T_0} + (r - c)\tau_{(0)}^{\text{on}} > x, B = 0] \\ &= \frac{\mathbb{E}\tau^{\text{off}}}{\mathbb{E}\tau^{\text{on}} + \mathbb{E}\tau^{\text{off}}} \mathbb{P}[Q_{T_0} + \tau_0^{\text{on}}(r - c) - c\tau_{(0)}^{\text{off}} > x] \\ &\quad + \frac{\mathbb{E}\tau^{\text{on}}}{\mathbb{E}\tau^{\text{on}} + \mathbb{E}\tau^{\text{off}}} \mathbb{P}[Q_{T_0} + (r - c)\tau_{(0)}^{\text{on}} > x]. \end{aligned} \tag{4.11}$$

(Note that $Q_0^{\text{P}} \equiv Q_{T_0}$.) Since Q_{T_0} and $\tau_{(0)}^{\text{on}}$ are independent, subexponential, and have asymptotically proportional tails, by applying Lemma A.3(ii) it follows that

$$\mathbb{P}[Q_{T_0} + (r - c)\tau_{(0)}^{\text{on}} > x] \sim \mathbb{P}[Q_{T_0} > x] + \mathbb{P}[(r - c)\tau_{(0)}^{\text{on}} > x] \quad \text{as } x \rightarrow \infty. \tag{4.12}$$

Next, $F_1 \in \mathcal{L}$ and Theorem 4.2 implies $\mathbb{P}[\tau_0^{\text{on}}(r - c) > x] = o(\mathbb{P}[Q_{T_0} > x])$ as $x \rightarrow \infty$, which in conjunction with Lemma A.3(i) and Lemma 4.1 yields

$$\mathbb{P}[Q_{T_0} + \tau_0^{\text{on}}(r - c) - c\tau_{(0)}^{\text{off}} > x] \sim \mathbb{P}[Q_{T_0} > x] \quad \text{as } x \rightarrow \infty. \tag{4.13}$$

Finally, by replacing asymptotic relations (4.12), (4.13), in (4.11), we obtain (4.8); combination of (4.6) and (4.8) gives (4.9). This completes the proof.

4.3. Multiplexing a long-tailed process with exponential processes

In this section we consider multiplexing one long-tailed on-off process with exponential processes (see Definition 4.2 below) in a fluid queue. The important conclusion to be drawn is that this queueing system is asymptotically interchangeable with a queueing system in which the on-off process is arriving alone and the exponential processes are replaced by their mean values. To reach this conclusion we need the following definition.

Definition 4.1. A distribution function F is *intermediate regular varying* ($F \in \mathcal{IR}$) if

$$\lim_{\delta \downarrow 1} \liminf_{t \rightarrow \infty} \frac{\bar{F}(\delta t)}{\bar{F}(t)} = 1.$$

Remark. For recent results on distributions of intermediate regular variation we refer the reader to [19]. Some basic properties of \mathcal{IR} are: $\mathcal{IR} \subset \mathcal{S}$; $\mathcal{R} \subset \mathcal{IR}$. Also, it is not very difficult to see that $\mathcal{IR} \subset \mathcal{S}_d$. Therefore, all of the results that we have obtained up to now apply for \mathcal{IR} . In addition, directly from the definition it can be shown that if $F \in \mathcal{IR}$ and $\int_0^\infty \bar{F}(t) dt < \infty$, then $F_1 \in \mathcal{IR}$.

Under the general large deviation Gärtner–Ellis conditions (see [52]) on the arrival process, it can be proved that the queue length distribution is exponentially bounded. To avoid stating Gärtner–Ellis conditions, we will define an arrival process e_t to be *exponential*, if whenever this process is fed into a constant server fluid queue, the queue length distribution is exponentially bounded.

Definition 4.2. We say that a stationary and ergodic arrival process e_t is *exponential* if for any server capacity $c > \mathbb{E}e_t$ there exist $K \equiv K(c)$ and $\delta \equiv \delta(c) > 0$ such that

$$\mathbb{P} \left[\sup_{t \geq 0} \int_{-t}^0 (e_u - c) du > x \right] \leq K e^{-\delta x}.$$

Remark. The main examples when the conditions of this definition are satisfied, i.e. Gärtner–Ellis conditions hold, are finite state space Markov chains or processes. Also, in terms of the on-off processes the conditions will hold whenever the distribution of on periods is exponentially bounded and off periods have a finite mean.

Recall that F and F_1 represent the distribution and the integrated tail distribution of an on period, respectively.

Theorem 4.4. Consider a single server queue with a capacity c and two independent arrival streams e_t and a_t . Assume that e_t is an exponential process (as in Definition 4.2) and a_t is an on-off process with rate r , $F \in \mathcal{LR}$, and generally distributed off periods with a finite mean. If $\mathbb{E}(e_t + a_t) < c$, $r > c' \stackrel{\text{def}}{=} c - \mathbb{E}e_t$, then the queue asymptotics of this queueing system is equal to the queue asymptotics in which only the on-off process arrives and the server capacity is replaced by c' , i.e. it is given by equation (4.9) in which c is replaced by c' .

Remark. (i) In [12, 13], a precise asymptotics of the embedded queue distribution was obtained for multiplexing on-off sources, one of which had regularly varying on periods while the others had exponentially distributed on periods. A similar setting with intermediately varying on periods was investigated in [49]. However, in both papers the equivalence relation of the original system to the system in which the exponential process is replaced by its mean has not been observed. (ii) The assumption of e_t being exponential can be weakened to $\mathbb{P}[\sup_{t \geq 0} \int_{-t}^0 (e_u - c) du > x] = o(F_1(x))$, for all $c > \mathbb{E}e_t$, without any changes in the proof.

Proof. Upper bound. For any $\epsilon > 0$, we can make the following decomposition

$$\begin{aligned} W_t &= \int_{-t}^0 (e_u + a_u - c) du = \int_{-t}^0 (e_u - \mathbb{E}e_0 - \epsilon) du + \int_{-t}^0 (a_u - (c - \mathbb{E}e_0 - \epsilon)) du \\ &\stackrel{\text{def}}{=} W_t^e + W_t^s. \end{aligned}$$

Call $c_\epsilon \stackrel{\text{def}}{=} c - \mathbb{E}e_0 - \epsilon$. Observe that for all sufficiently small $\epsilon > 0$, such that $\mathbb{E}(e_t + a_t) + \epsilon < c$, $\mathbb{P}[\sup_{t \geq 0} W_t^s \leq x]$ represents the queue length distribution in a stable on-off queue with arrival process a_t , and service capacity c_ϵ . This implies that $\mathbb{P}[\sup_{t \geq 0} W_t^s \leq x] \in \mathcal{LR}$ (by using Theorem 4.3).

Further,

$$\begin{aligned} \mathbb{P}[Q_t > x] &= \mathbb{P}\left[\sup_{t \geq 0} (W_t^e + W_t^s) > x\right] \\ &\leq \mathbb{P}\left[\sup_{t \geq 0} W_t^e + \sup_{t \geq 0} W_t^s > x\right] \\ &\sim \mathbb{P}\left[\sup_{t \geq 0} W_t^s > x\right] \quad \text{as } x \rightarrow \infty, \end{aligned} \tag{4.14}$$

where (4.14) follows from Lemma A.3(i), since $\sup_{t \geq 0} W_t^e$ and $\sup_{t \geq 0} W_t^s$ are two independent random variables, with $\mathbb{P}[\sup_{t \geq 0} W_t^e > x] = o(\mathbb{P}[\sup_{t \geq 0} W_t^s > x])$ as $x \rightarrow \infty$. Now,

$$\lim_{x \rightarrow \infty} \frac{\mathbb{P}[\sup_{t \geq 0} W_t^s > x]}{\int_{x/(r-c_\epsilon)}^\infty \mathbb{P}[\tau^{\text{on}} > u] du} = K_\epsilon, \tag{4.15}$$

where K_ϵ is given by equation (4.10) in Theorem 4.3, with c_ϵ in place of c . Consequently, this leads to

$$\limsup_{x \rightarrow \infty} \frac{\mathbb{P}[Q_t > x]}{\int_{x/(r-c')}^\infty \mathbb{P}[\tau^{\text{on}} > u] du} \leq K_\epsilon \limsup_{x \rightarrow \infty} \frac{\int_{x/(r-c_\epsilon)}^\infty \mathbb{P}[\tau^{\text{on}} > u] du}{\int_{x/(r-c')}^\infty \mathbb{P}[\tau^{\text{on}} > u] du}. \tag{4.16}$$

Finally, if we let $\epsilon \rightarrow 0$ in (4.16), we obtain

$$\limsup_{x \rightarrow \infty} \frac{\mathbb{P}[Q_t > x]}{\int_{x/(r-c')}^\infty \mathbb{P}[\tau^{\text{on}} > u] du} \leq K', \tag{4.17}$$

where K' , similarly to K_ϵ , is given by equation (4.10) in Theorem 4.3, with c' in place of c . This concludes the proof of the upper bounds.

Lower bound. For the lower bound we consider a different decomposition of W_t , i.e., we redefine W_t^e and W_t^s as follows ($\epsilon > 0$)

$$\begin{aligned} W_t &= \int_{-t}^0 (e_u + a_u - c) du = - \int_{-t}^0 (\mathbb{E}e_0 - \epsilon - e_u) du + \int_{-t}^0 (a_u - (c - \mathbb{E}e_0 + \epsilon)) du \\ &\stackrel{\text{def}}{=} -W_t^e + W_t^s. \end{aligned}$$

Also, redefine $c_\epsilon \stackrel{\text{def}}{=} c - \mathbb{E}e_0 + \epsilon$. As in the upper bound case, for $\epsilon < \mathbb{E}e_0 + \mathbb{E}a_0 - c$, $\mathbb{P}[\sup_{t \geq 0} W_t^s \leq x]$ represents a queue length distribution in a stable on-off queue with arrival process a_t , and service capacity c_ϵ . Hence, $\mathbb{P}[\sup_{t \geq 0} W_t^s \leq x] \in \mathcal{LR}$. Further,

$$\begin{aligned} \mathbb{P}[Q_t > x] &= \mathbb{P}\left[\sup_{t \geq 0} (-W_t^e + W_t^s) > x\right] \\ &\geq \mathbb{P}\left[\inf_{t \geq 0} -W_t^e + \sup_{t \geq 0} W_t^s > x\right] \end{aligned} \tag{4.18}$$

$$\begin{aligned} &= \mathbb{P}\left[-\sup_{t \geq 0} W_t^e + \sup_{t \geq 0} W_t^s > x\right] \\ &\sim \mathbb{P}\left[\sup_{t \geq 0} W_t^s > x\right] \quad \text{as } x \rightarrow \infty, \end{aligned} \tag{4.19}$$

where in (4.18) we use the fact that $\sup(f + g) \geq \inf f + \sup g$ for any two functions $f, g : \mathbb{R} \rightarrow \mathbb{R}$; asymptotics in (4.19) follows from the independence of a_t and e_t , $\mathbb{P}[\sup_{t \geq 0} W_t^s \leq x] \in \mathcal{L}$, $\mathbb{P}[\sup_{t \geq 0} W_t^e < \infty] = 1$, and Lemma 4.1(ii). Consequently, (4.19) leads to

$$\liminf_{x \rightarrow \infty} \frac{\mathbb{P}[Q_t > x]}{\int_{x/(r-c')}^{\infty} \mathbb{P}[\tau^{\text{on}} > u] du} \geq K_\epsilon \liminf_{x \rightarrow \infty} \frac{\int_{x/(r-c_\epsilon)}^{\infty} \mathbb{P}[\tau^{\text{on}} > u] du}{\int_{x/(r-c')}^{\infty} \mathbb{P}[\tau^{\text{on}} > u] du}, \tag{4.20}$$

where again K_ϵ is computed from Theorem 4.3 with c_ϵ in place of c . Consequently, using the same arguments as in (4.17) we arrive at

$$\liminf_{x \rightarrow \infty} \frac{\mathbb{P}[Q_t > x]}{\int_{x/(r-c')}^{\infty} \mathbb{P}[\tau^{\text{on}} > u] du} \geq K', \tag{4.21}$$

which together with (4.17) concludes the proof.

4.4. Subexponential $M/G/\infty$ arrival process

In the next theorem we obtain a tight lower bound for the fluid queue asymptotics with $M/G/\infty$ arrivals. For this fluid queue we denote its queue content process by Q_t^∞ .

4.4.1. Lower bound

Theorem 4.5. *Let $\rho \stackrel{\text{def}}{=} \mathbb{E}A_t^{\infty,s} = \Lambda r \mathbb{E}\tau^{\text{on}} < c$. If $r + \rho > c$, and $\tau^{\text{on}} \in \mathcal{LR}$, then*

$$\liminf_{x \rightarrow \infty} \frac{\mathbb{P}[Q_t^\infty > x]}{\int_{x/(r+\rho-c)}^{\infty} \mathbb{P}[\tau^{\text{on}} > u] du} \geq \frac{\Lambda r}{c - \rho}.$$

Proof. Let $\underline{A}^{\infty,y} = \{A_t^{\infty,y}, t \geq 0\}$, $\overline{A}^{\infty,y} = \{\overline{A}_t^{\infty,y}, t \geq 0\}$, $y > 0$, be two independent $M/G/\infty$ type processes with Poisson arrival rate and with on distributions respectively given as $\underline{A}^y = \mathbb{P}[\tau^{\text{on}} \leq y] \Lambda$, $\overline{A}^y = \mathbb{P}[\tau^{\text{on}} > y] \Lambda$, $F_{\underline{y}}(x) = \mathbb{P}[\tau^{\text{on}} \leq x] / \mathbb{P}[\tau^{\text{on}} \leq y]$, $0 \leq x \leq y$, $F_{\overline{y}}(x) = \mathbb{P}[y < \tau^{\text{on}} \leq x] / \mathbb{P}[\tau^{\text{on}} > y]$, $x > y$; all three processes are assumed to have the same parameter r . Now, we claim that

$$A^\infty \stackrel{\text{d}}{=} \underline{A}^{\infty,y} + \overline{A}^{\infty,y}, \tag{4.22}$$

where $\stackrel{\text{d}}{=}$ stands for equality in distribution. Observe that, for a fixed parameter r , any $M/G/\infty$ process is uniquely defined with Poisson arrival times $\{T_n\}$, and the lengths of on periods $\{\tau_n^{\text{on}}\}$; likewise, the pair $\{T_n\}, \{\tau_n^{\text{on}}\}$ uniquely defines a compound Poisson process (see [17, p. 90]), that is a piecewise constant process with Poisson jump times $\{T_n\}$, and jump sizes $\{\tau_n^{\text{on}}\}$. Hence, proving (4.22) is equivalent to

$$Z \stackrel{\text{d}}{=} \underline{Z}^y + \overline{Z}^y, \tag{4.23}$$

where $Z, \underline{Z}^y, \overline{Z}^y$, are the compound Poisson processes corresponding to $A^\infty, \underline{A}^{\infty,y}, \overline{A}^{\infty,y}$, respectively. Since $Z, \underline{Z}^y, \overline{Z}^y$, are processes with stationary independent increments, (4.23) is equivalent to

$$Z_t \stackrel{\text{d}}{=} \underline{Z}_t^y + \overline{Z}_t^y, \quad t \geq 0. \tag{4.24}$$

Evidently (4.24) is implied by

$$\begin{aligned} & \mathbb{E} \exp\{-s(\underline{Z}_t^y + \overline{Z}_t^y)\} \\ &= \mathbb{E} \exp\{-s\underline{Z}_t^y\} \mathbb{E} \exp\{-s\overline{Z}_t^y\} \\ &= \exp\left(-t\underline{\Lambda}^y \int_0^\infty (1 - e^{-su}) dF_{\underline{Y}}(u)\right) \exp\left(-t\overline{\Lambda}^y \int_0^\infty (1 - e^{-su}) dF_{\overline{Y}}(u)\right) \end{aligned} \quad (4.25)$$

$$\begin{aligned} &= \exp\left(-t\Lambda \int_0^y (1 - e^{-su}) dF(u)\right) \exp\left(-t\Lambda \int_y^\infty (1 - e^{-su}) dF(u)\right) \\ &= \exp\left(-t\Lambda \int_0^\infty (1 - e^{-su}) dF(u)\right) = \mathbb{E} \exp\{-sZ_t\}; \end{aligned} \quad (4.26)$$

equalities (4.25), (4.26), are well known expressions for compound Poisson processes (see [17], p. 94). This proves (4.22).

Note that process $\underline{A}_t^{\infty,y}$ has bounded on periods, and in conclusion, it is an *exponential* process (i.e., it satisfies Definition 4.2). Also, $\overline{A}_t^{\infty,y} \geq r\mathbf{1}(\overline{A}_t^{\infty,y} > 0)$. Therefore, $\mathbb{P}[Q_t^\infty > x]$ is stochastically larger than a queueing process \underline{Q}_t^∞ obtained by feeding $\underline{A}_t^{\infty,y} + r\mathbf{1}(\overline{A}_t^{\infty,y} > 0)$ into it. Let $Q_t^{a_y}$ be a queueing process with a subexponential on-off arrival process $a_t^y \stackrel{\text{def}}{=} r\mathbf{1}(\overline{A}_t^{\infty,y} > 0)$, and a server capacity $c_y = c - \mathbb{E}\underline{A}_t^{\infty,y} = c - r\mathbb{P}[\tau^{\text{on}} \leq y]\Lambda\mathbb{E}\tau^{\text{on}}$. Here, by Theorem 4.4, we obtain

$$\mathbb{P}[Q_t^\infty > x] \geq \mathbb{P}[\underline{Q}_t^\infty > x] \sim \mathbb{P}[Q_t^{a_y} > x] \quad \text{as } x \rightarrow \infty. \quad (4.27)$$

In addition, a_t^y has off period $I^{\text{off},y}$ exponentially distributed with parameter $\overline{\Lambda}^y$, on period $I^{\text{on},y}$ with mean $\mathbb{E}I^{\text{on},y} = 1/\overline{\Lambda}^y (\exp\{\overline{\Lambda}^y \mathbb{E}\tau^{\text{on},y}\} - 1) = 1/\overline{\Lambda}^y (\exp\{\Lambda\mathbb{E}\tau^{\text{on}}\mathbf{1}(\tau^{\text{on}} > y)\} - 1)$, and asymptotics (by Theorem 3.5)

$$\begin{aligned} \mathbb{P}[I^{\text{on},y} > t] &\sim \exp\{\overline{\Lambda}^y \mathbb{E}\tau^{\text{on},y}\} \mathbb{P}[\tau^{\text{on},y} > t] \quad \text{as } t \rightarrow \infty \\ &= \frac{\exp\{\Lambda\mathbb{E}\tau^{\text{on}}\mathbf{1}(\tau^{\text{on}} > y)\}}{\mathbb{P}[\tau^{\text{on}} > y]} \mathbb{P}[\tau^{\text{on}} > t]. \end{aligned} \quad (4.28)$$

Hence, Theorem 4.3 and (4.28) lead to

$$\mathbb{P}[Q_t^{a_y} > x] \sim K_y \int_{x/(r-c_y)}^\infty \mathbb{P}[\tau^{\text{on}} > u] du \quad \text{as } x \rightarrow \infty, \quad (4.29)$$

where

$$K_y \stackrel{\text{def}}{=} \left(\frac{r - c_y}{c_y \mathbb{E}I^{\text{off},y} - (r - c_y) \mathbb{E}I^{\text{on},y}} + \frac{1}{\mathbb{E}I^{\text{off},y} + \mathbb{E}I^{\text{on},y}} \right) \frac{\exp\{\Lambda\mathbb{E}\tau^{\text{on}}\mathbf{1}(\tau^{\text{on}} > y)\}}{\mathbb{P}[\tau^{\text{on}} > y]}.$$

Combining (4.27) and (4.29) produces

$$\liminf_{x \rightarrow \infty} \frac{\mathbb{P}[Q_t^\infty > x]}{\int_{x/(r-c')}^\infty \mathbb{P}[\tau^{\text{on}} > u] du} \geq K_y \liminf_{x \rightarrow \infty} \frac{\int_{x/(r-c_y)}^\infty \mathbb{P}[\tau^{\text{on}} > u] du}{\int_{x/(r-c')}^\infty \mathbb{P}[\tau^{\text{on}} > u] du}, \quad (4.30)$$

where $c' \stackrel{\text{def}}{=} c - \rho$. Here, observe that $K_y \rightarrow \Lambda r/c'$ as $y \rightarrow \infty$, and, by assumption $F \in \mathcal{LR}$,

$$\lim_{y \rightarrow \infty} \liminf_{x \rightarrow \infty} \frac{\int_{x/(r-c_y)}^\infty \mathbb{P}[\tau^{\text{on}} > u] du}{\int_{x/(r-c')}^\infty \mathbb{P}[\tau^{\text{on}} > u] du} = 1.$$

Finally, if we take the limit with respect to y in (4.30), we obtain the statement of the theorem.

4.4.2. *Precise queue asymptotics.* Let $Q_n^{P,\infty}$ be the queue size observed at the beginning of the n th activity period in the $M/G/\infty$ arrival process.

Theorem 4.6. *Let $\rho = \mathbb{E}A_t^{s,\infty} = \Lambda r \mathbb{E}\tau^{\text{on}} < c$. If $c \leq r$, and τ^{on} is regularly varying with non-integer exponent $\alpha > 1$, then*

$$\lim_{x \rightarrow \infty} \frac{\mathbb{P}[Q_t^{P,\infty} > x]}{\int_{x/(\rho+r-c)}^{\infty} \mathbb{P}[\tau^{\text{on}} > u] du} = \Lambda \left(\frac{r}{c - \rho} - 1 \right).$$

Proof. Proof follows directly from Theorem 4.1 and Theorem 3.6, by taking $A_n \stackrel{\text{def}}{=} D_{c,n}$, $C_n = cI_n^{\text{off}}$, and observing that $\mathbb{E}(A_n - C_n) = e^{\Lambda \mathbb{E}\tau^{\text{on}}} (\Lambda r \mathbb{E}\tau^{\text{on}} - c) / \Lambda$.

Theorem 4.7. *Let $\rho = \mathbb{E}A_t^{s,\infty} = \Lambda r \mathbb{E}\tau^{\text{on}} < c$. If Conjecture 3.1 holds, $c \leq r$, and τ^{on} is regularly varying with non-integer exponent $\alpha > 1$, then*

$$\lim_{x \rightarrow \infty} \frac{\mathbb{P}[Q_t^{\infty} > x]}{\int_{x/(\rho+r-c)}^{\infty} \mathbb{P}[\tau^{\text{on}} > u] du} = \Lambda \frac{r}{c - \rho}.$$

Proof. This theorem follows from Theorem 4.8, Conjecture 3.1, and exactly the same arguments as in the proof of Theorem 4.3. We omit the details.

Remark. The asymptotic result in this theorem is the same as the lower bound obtained in Theorem 4.5.

4.5. Finite number of subexponential on-off processes

For a finite number of long-tailed on-off processes we can easily obtain the upper bound given by Theorem C.1. Similarly, utilization of $A_t^N \geq r \mathbf{1}(A_t^N > 0)$, Theorem 3.3 and Theorem 4.3 can easily produce a lower bound. Unfortunately, these bounds are very weak and for this reason we resort back to Theorems 4.5, 4.6, and 4.7. Based on these results in [31, 32, 33] we have suggested an approximation for the finite number of on-off processes. In the same papers this approximation was tested using simulation experiments.

5. Conclusion

In this paper we have established a precise asymptotic characterization of the activity period of an arrival process obtained by multiplexing on-off processes with exponential off periods and subexponential on periods. This characterization has been done both for a finite number of processes as well as for the limiting $M/G/\infty$ case.

For a simple subexponential on-off fluid flow queue we have obtained a precise asymptotic relation between the Palm queue distribution and the time average queue distribution. Furthermore, exponential processes, when multiplexed with a subexponential on-off process, have been shown to contribute to the large buffer asymptotics only through their mean value.

In the limiting $M/G/\infty$ case (e.g. large number of subexponential on-off processes) with regularly varying on periods with non-integer exponents we have obtained a precise queue asymptotics observed at the beginning of the arrival process activity periods. The asymptotic

time average queue lower bound has been derived under more general assumptions of intermediately varying on periods.

Based on these asymptotic results, a computationally efficient approximation was suggested in [31, 32, 33] for the large buffer probabilities of finitely many subexponential on-off processes. The accuracy of this approximation was verified using extensive simulation experiments.

The results in this paper bring us closer to understanding the subexponential queueing asymptotics of multiplexed long-tailed processes. The precision and low computational complexity of the $M/G/\infty$ approximation has a practical impact on improving the efficacy of ATM admission controllers.

Appendix A. Basic results on subexponential and long-tailed distributions

In what follows we will state a few important results from the literature on subexponential distributions. The general relation between \mathcal{S} and \mathcal{L} is the following.

Lemma A.1. [8] $\mathcal{S} \subset \mathcal{L}$.

Lemma A.2. If $F \in \mathcal{L}$ then $(1 - F(x))e^{\alpha x} \rightarrow \infty$ as $x \rightarrow \infty$, for all $\alpha > 0$.

Note. Lemma A.2 clearly shows that for long-tailed distributions Cramér type conditions are not satisfied.

The proof of the following result can be found in [26].

Lemma A.3. Let $F \in \mathcal{S}$. Then,

- (i) If G is a probability distribution such that $\bar{G}(x) = o(\bar{F}(x))$ as $x \rightarrow \infty$, then $\overline{F * G}(x) \sim \bar{F}(x)$ as $x \rightarrow \infty$.
- (ii) If $\lim_{x \rightarrow \infty} \bar{G}(x)/\bar{F}(x) = c \in (0, \infty)$, where G is a distribution function on $[0, \infty)$, then $G \in \mathcal{S}$, and $\overline{F * G}(x) \sim \bar{F}(x) + \bar{G}(x)$ as $x \rightarrow \infty$.

The next result is due to Athreya and Ney (see [8, pp. 147–150]).

Lemma A.4. If $F \in \mathcal{S}$, then

- (i) $\overline{F^{*n}}(x)/\bar{F}(x) \rightarrow n$ as $x \rightarrow \infty$, for all $n \in \mathbb{N}$.
- (ii) For each $\epsilon > 0$ there exists a constant $C_\epsilon (< \infty)$ such that $\overline{F^{*n}}(x) \leq C_\epsilon(1 + \epsilon)^n \bar{F}(x)$ for all x and n .

This lemma directly gives the asymptotics of a renewal measure with the following Pollaczek–Khintchine representation

$$G(x) \stackrel{\text{def}}{=} \sum_{n=0}^{\infty} \gamma^n \overline{F^{*n}}(x). \tag{A.1}$$

Using dominated convergence and the previous lemma it is easy to prove the following very useful theorem.

Theorem A.1. If $F \in \mathcal{S}$, and $-1 < \gamma < 1$, then

$$\lim_{x \rightarrow \infty} \frac{\bar{G}(x)}{\bar{F}(x)} = \frac{\gamma}{(1 - \gamma)^2}.$$

Proof. Follows easily from Lemma A.4 and dominated convergence (see [6]).

Sometimes when F is absolutely continuous, i.e. F has a density function f , it is of interest to calculate the density g of G . If we take a derivative in (A.1) with respect to x we obtain

$$g(x) \stackrel{\text{def}}{=} \sum_{n=0}^{\infty} \gamma^n f^{\otimes n}(x), \tag{A.2}$$

where $f^{\otimes n}(x)$ is the density of F^{n*} ; $f^{\otimes 2}(x) = \int_0^x f(x-u)f(u) du$, and $f^{\otimes(n+1)}(x) = \int_0^x f^{\otimes n}(x-u)f(u) du$. The investigation of the asymptotics of g requires the investigation of the asymptotics of $f^{\otimes n}(x)$. This motivated the introduction of subexponential density functions in [40].

Definition A.1. A function $f : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ such that $f(x) > 0$ on $[A, \infty)$ for some $A \in \mathbb{R}_+$ belongs to the class \mathcal{S}_d if f is long-tailed, and

$$\lim_{x \rightarrow \infty} \frac{f^{\otimes 2}(x)}{f(x)} = 2d.$$

Then, equivalent results to Lemma A.4, and Theorem A.1 (in its most general form) were obtained in Theorem 3.2 of [40]. For convenience we state here simplified versions of these results.

Lemma A.5. Let $f \in \mathcal{S}_d$, and $\int_0^\infty f(x) dx = 1$. Then, both (i) and (ii) of Lemma A.4 are true with \bar{F} replaced by f and $\overline{F^{*n}}(x)$ replaced by $f^{\otimes n}$.

Theorem A.2. If $f \in \mathcal{S}_d$, with $\int_0^\infty f(x) dx = 1$, and $-1 < \gamma < 1$, then

$$g(x) \sim \frac{\gamma}{(1-\gamma)^2} f(x) \quad \text{as } x \rightarrow \infty.$$

An extensive treatment of subexponential distributions and further references can be found in [18, 39]. For a recent survey of the application of subexponential distributions in queueing theory the reader is referred to [5].

Appendix B. Proof of Theorem 3.6

As we have already mentioned, the proof of this result is based on Karamata’s Tauberian/-Abelian theorem for distribution functions of regular variation. This theorem relates the tail behavior of a distribution function to the asymptotic behavior of its Laplace transform at the origin. For convenience we state the following result due to Bingham and Doney ([10], [11, p. 333]).

Let F be a distribution function on $[0, \infty)$, and let $\tilde{F}(s)$ be its Laplace–Stieltjes transform. Denote by $m_n = \mathbb{E}X^n = \int_{[0, \infty)} x^n dF(x)$, $n = 0, 1, \dots$. When $m_n < \infty$, $\tilde{F}(s)$ may be expanded in a Taylor series as far as the s^n term:

$$\tilde{F}(s) = \sum_{k=0}^n m_k (-s)^k / k! + o(s^n) \quad \text{as } s \downarrow 0.$$

To compare the tail behavior of F with the behavior of \tilde{F} at the origin, one needs to eliminate the Taylor polynomial $\sum_{k=0}^n m_k(-s)^k/k!$. This may be done by subtraction or repeated differentiation, i.e. let

$$f_n(s) \stackrel{\text{def}}{=} (-1)^{n+1} \left\{ \tilde{F}(s) - \sum_{k=0}^n m_k(-s)^k/k! \right\},$$

$$g_n(s) \stackrel{\text{def}}{=} d^n f_n(s)/ds^n = m_n - (-1)^n \tilde{F}^{(n)}(s), \quad n \geq 0,$$

where $\tilde{F}^{(n)}(s)$ denotes the n th derivative of $\tilde{F}(s)$; in general for any function G we will use $G^{(n)}$ to denote its n th derivative.

Theorem B.1. *Let l be a slowly varying function, $n \in \mathbb{N}_0$, and $\alpha = n + \beta$, $0 < \beta < 1$ (α non-integer). Then, the following are equivalent*

$$f_n(s) \sim s^\alpha l(1/s) \quad \text{as } s \downarrow 0, \tag{B.1}$$

$$g_n(s) \sim \frac{\Gamma(\alpha + 1)}{\Gamma(\beta + 1)} s^\beta l(1/s) \quad \text{as } s \downarrow 0, \tag{B.2}$$

$$(-1)^{n+1} \tilde{F}^{(n+1)}(s) \sim \frac{\Gamma(\alpha + 1)}{\Gamma(\beta)} s^{\beta-1} l(1/s) \quad \text{as } s \downarrow 0, \tag{B.3}$$

$$1 - F(x) \sim \frac{(-1)^n}{\Gamma(1 - \alpha)} x^{-\alpha} l(x) \quad \text{as } x \rightarrow \infty, \tag{B.4}$$

where Γ stands for the gamma function.

Now, in order to determine the tail behavior of $D_{c,n}$, we will investigate asymptotics of its LS transform at the origin. Note that without loss of generality we can set $r = 1$, since $D_{c,r} = rD_{c/r,1}$; $D_{c,r}$ stands for D_c where the arrival process has parameter r . To simplify the notation we use $D \equiv D_c$, $\tau \equiv \tau^{\text{on}}$. Slight modification of Theorem 2.10, [21], reads as follows.

Theorem B.2. *Let $0 \leq c \leq 1$ ($r = 1$), $\mathbb{E}\tau < \infty$. Then, for $s > 0$,*

$$\frac{c - \Lambda \mathbb{E}\tau \tilde{F}_1(s)}{c - \Lambda \mathbb{E}D \tilde{D}_1(s)} = 1 - \Lambda \int_0^\infty \mathbb{E}[\exp\{-s(\tau - ct)\} \mathbf{1}(\tau \geq t)] \exp\{-f(s, t)\} dt, \tag{B.5}$$

where

$$f(s, t) = \Lambda \{t(1 - \tilde{F}(s)) + \mathbb{E}[\tau \exp(-s\tau)] - \mathbb{E}[(\tau - t) \exp(-s\tau) \mathbf{1}(\tau \geq t)]\},$$

$\tilde{F}(s) = \mathbb{E}e^{-s\tau}$, and $\tilde{F}_1(s) = (1 - \mathbb{E}e^{-s\tau})/(s\mathbb{E}\tau)$, $\tilde{D}_1(s) = (1 - \mathbb{E}e^{-sD})/(s\mathbb{E}D)$, are the Laplace–Stieltjes transforms of the integrated tail distributions of F and D , respectively; $\mathbb{E}D = (\mathbb{E}\tau - c/\Lambda) e^{\Lambda \mathbb{E}\tau} + c/\Lambda$.

Proof of Theorem 3.6. Let $r = 1$. In order to simplify the usage of Theorem B.1, without loss of generality we assume that $1 - F(x) \sim (-1)^n x^{-\alpha} l(x) / \Gamma(1 - \alpha)$, $\alpha = n + \beta$, $n \in \mathbb{N}$, $0 < \beta < 1$. Let us write the integral in (B.5) as

$$\phi(s) = \int_0^\infty g(s, t) e^{-f(s,t)} dt,$$

where $g(s, t) \stackrel{\text{def}}{=} \mathbb{E}\{e^{-s(\tau-ct)} \mathbf{1}(\tau \geq t)\}$. Next, we want to find the n th derivative of $\phi(s)$ and compare it with the n th derivative of the left-hand side of (B.5) as $s \downarrow 0$. The *main technical difficulty* in doing this is to prove the following two lemmas.

Lemma B.1. *For any $\rho > \epsilon > 0$, there exist $s_0 > 0, t_0 > 0$, such that for all $0 < s < s_0, t > t_0, 0 \leq k \leq n$,*

$$(1 - \epsilon)\rho^k t^k e^{-s(\rho+\epsilon)t-\rho} \leq (-1)^k \frac{\partial^k}{\partial s^k} e^{-f(s,t)} \leq (1 + \epsilon)\rho^k t^k e^{-s(\rho-\epsilon)t-\rho}. \tag{B.6}$$

Proof. Given at the end of this section.

Lemma B.1 is crucial in proving the following result which directly leads to the proof of Theorem 3.6. Unfortunately, the proof of this lemma is very technical and, due to the space limitations, we do not present it here.

Lemma B.2. *As $s \downarrow 0$*

$$\phi^{(n)}(s) \sim \frac{d^n}{ds^n} \int_0^\infty \mathbb{E}\{e^{-s(\tau-ct)} \mathbf{1}(\tau \geq t)\} e^{-\rho st-\rho} dt, \tag{B.7}$$

where $\rho = \Lambda \mathbb{E}\tau$.

Proof. Given in [31].

At this point we are ready to finish the proof of Theorem 3.6. By using Lemma B.2 we derive

$$\begin{aligned} & \lim_{s \downarrow 0} s^{1-\beta} (l(1/s))^{-1} \phi^{(n)}(s) \\ &= \lim_{s \downarrow 0} s^{1-\beta} (l(1/s))^{-1} e^{-\rho} \frac{d^n}{ds^n} \mathbb{E} e^{-s\tau} \int_0^\tau e^{-(\rho-c)st} dt \\ &= \lim_{s \downarrow 0} s^{1-\beta} (l(1/s))^{-1} e^{-\rho} \frac{d^n}{ds^n} \mathbb{E} e^{-s\tau} \frac{e^{-(\rho-c)s\tau} - 1}{-s(\rho - c)} \\ &= \lim_{s \downarrow 0} s^{1-\beta} (l(1/s))^{-1} e^{-\rho} \frac{d^n}{ds^n} \mathbb{E} \frac{e^{-s\tau} - e^{-(\rho+1-c)s\tau}}{s(\rho - c)} \\ &= \frac{e^{-\rho} \mathbb{E}\tau}{\rho - c} \lim_{s \downarrow 0} s^{1-\beta} (l(1/s))^{-1} \frac{d^n}{ds^n} \left[\mathbb{E} \frac{1 - e^{-(\rho+1-c)s\tau}}{s\mathbb{E}\tau} - \mathbb{E} \frac{1 - e^{-s\tau}}{s\mathbb{E}\tau} \right] \\ &= \frac{e^{-\rho} \mathbb{E}\tau}{\rho - c} \lim_{s \downarrow 0} s^{1-\beta} (l(1/s))^{-1} \frac{d^n}{ds^n} [(\rho + 1 - c)\tilde{F}_1(s(\rho + 1 - c)) - \tilde{F}_1(s)]; \tag{B.8} \end{aligned}$$

recall that $\tilde{F}_1(s)$ is the LS transform of the integrated tail distribution of τ . From the assumption on F it follows that

$$\begin{aligned} 1 - F_1(x) &\sim (-1)^n x^{-(\alpha-1)} l(x) / (\Gamma(1 - \alpha)(\alpha - 1)\mathbb{E}\tau) \\ &= (-1)^{n-1} x^{-(\alpha-1)} l(x) / (\Gamma(2 - \alpha)\mathbb{E}\tau) \quad x \rightarrow \infty. \end{aligned}$$

Consequently, by Theorem B.1, equation (B.3),

$$\tilde{F}_1^{(n)}(s) \sim (-1)^n \frac{\Gamma(\alpha)}{\Gamma(\beta)\mathbb{E}\tau} s^{\beta-1} l(1/s), \quad \text{as } s \downarrow \infty. \tag{B.9}$$

The latter result, when replaced in (B.8), produces

$$-\Lambda \lim_{s \downarrow 0} s^{1-\beta} l(1/s)^{-1} \phi^{(n)}(s) = \frac{\Lambda e^{-\rho}}{\rho - c} \times \frac{(-1)^n \Gamma(\alpha)}{\Gamma(\beta)} [1 - (\rho + 1 - c)^\alpha], \quad (\text{B.10})$$

which represents the result of taking the operator $\lim_{s \downarrow 0} s^{1-\beta} l(1/s)^{-1} d^n/ds^n$ on the right-hand side of (B.5). To finish the proof we have to compute the result of applying the same operator on the left-hand side of (B.5). Let us start with the derivative

$$\frac{d^n}{ds^n} \frac{c - \Lambda \mathbb{E} \tau \tilde{F}_1(s)}{c - \Lambda \mathbb{E} D \tilde{D}_1(s)} = \frac{-\Lambda \mathbb{E} \tau \tilde{F}_1^{(n)}(s)}{c - \Lambda \mathbb{E} D \tilde{D}_1(s)} + \frac{(c - \Lambda \mathbb{E} \tau \tilde{F}_1(s)) \Lambda \mathbb{E} D \tilde{D}_1^{(n)}(s)}{(c - \Lambda \mathbb{E} D \tilde{D}_1(s))^2} + R_{n-1}(s), \quad (\text{B.11})$$

where $R_{n-1}(s)$ is a rational function that contains only the first $n - 1$ derivatives of $\tilde{D}_1(s)$, and $\tilde{F}_1(s)$. Note that $|\tilde{F}_1^{(k)}(s)| \leq |\tilde{F}_1^{(k)}(0)| < \infty$, for $0 \leq k \leq n - 1$, $s \geq 0$. Furthermore, by induction and by taking successive derivatives in (B.5), it is easy to show that $|\tilde{D}_1^{(k)}(s)| \leq |\tilde{D}_1^{(k)}(0)| < \infty$, for $0 \leq k \leq n - 1$, $s \geq 0$. Therefore,

$$\lim_{s \downarrow 0} |R_{n-1}(s)| < \infty. \quad (\text{B.12})$$

Combining (B.9) and $\mathbb{E} D = (\mathbb{E} \tau - c/\Lambda) e^\rho + c/\Lambda$ we obtain

$$\frac{-\Lambda \mathbb{E} \tau \tilde{F}_1^{(n)}(s)}{c - \Lambda \mathbb{E} D \tilde{D}_1(s)} \sim \frac{\Lambda e^{-\rho}}{\rho - c} \times \frac{(-1)^n \Gamma(\alpha)}{\Gamma(\beta)} s^{\beta-1} l(1/s) \quad s \rightarrow 0. \quad (\text{B.13})$$

Finally, $\mathbb{E} D = (\mathbb{E} \tau - c/\Lambda) e^\rho + c/\Lambda$, (B.10), (B.11), (B.12), and (B.13), yield

$$\begin{aligned} \lim_{s \downarrow 0} s^{1-\beta} l(1/s) \tilde{D}_1^{(n)}(s) &= \frac{(c - \Lambda \mathbb{E} D)^2}{(c - \rho) \Lambda \mathbb{E} D} \times \frac{-\Lambda e^{-\rho}}{\rho - c} \times \frac{(-1)^n \Gamma(\alpha)}{\Gamma(\beta)} (\rho + 1 - c)^\alpha \\ &= e^\rho (\rho + 1 - c)^\alpha \frac{(-1)^n \Gamma(\alpha)}{\Gamma(\beta) \mathbb{E} D}, \end{aligned} \quad (\text{B.14})$$

which is, by Theorem B.1, equivalent to

$$\int_x^\infty \mathbb{P}[D > u] du \sim e^\rho (\rho + 1 - c)^\alpha (-1)^{n-1} x^{-(\alpha-1)} l(x) / \Gamma(2 - \alpha) \quad \text{as } x \rightarrow \infty;$$

finally, by the Monotone Density Theorem [11, p. 39]

$$\begin{aligned} \mathbb{P}[D > x] &\sim e^\rho (\rho + 1 - c)^\alpha (-1)^n x^{-\alpha} l(x) / \Gamma(1 - \alpha) \sim e^\rho (\rho + 1 - c)^\alpha \mathbb{P}[\tau > x] \\ &\sim e^\rho \mathbb{P}[\tau(\rho + 1 - c) > x] \end{aligned}$$

as $x \rightarrow \infty$; this finishes the proof of Theorem 3.6 for the case $r = 1$. For $r \neq 1$, $D_{c,r} = r D_{c/r,1}$; thus

$$\mathbb{P}[D_{c,r} > x] = \mathbb{P}[r D_{c/r,1} > x] \sim e^\rho (r + \rho - c)^\alpha \mathbb{P}[\tau > x] \quad \text{as } x \rightarrow \infty,$$

where $\rho = r \Lambda \mathbb{E} \tau$. This completes the proof of Theorem 3.6.

Proof of Lemma B.1. Let us start with the case $k = 0$. Since $\mathbb{E}[(\tau - t) \exp(-s\tau) \mathbf{1}(\tau \geq t)] \leq \mathbb{E}[(\tau - t) \mathbf{1}(\tau \geq t)] \rightarrow 0$, as $t \rightarrow \infty$, and $(1 - \tilde{F}(s))/s \rightarrow \mathbb{E}\tau$, $\mathbb{E}[\tau \exp(-s\tau)] \rightarrow \mathbb{E}\tau$ as $s \rightarrow 0$, it follows that for any $\epsilon > 0$, $\epsilon < \rho$, there exist $s_0 > 0$, $t_0 > 0$, such that for all $0 < s < s_0$, $t > t_0$

$$(1 - \epsilon) e^{-s(\rho+\epsilon)t-\rho} \leq e^{-f(s,t)} \leq (1 + \epsilon) e^{-s(\rho-\epsilon)t-\rho}. \tag{B.15}$$

For $1 \leq k \leq n$ let us first show that for any $\rho > \epsilon > 0$,

$$(1 - \epsilon) \rho^k t^k \leq (-1)^k e^{f(s,t)} \frac{\partial^k}{\partial s^k} e^{-f(s,t)} \leq (1 + \epsilon) \rho^k t^k, \tag{B.16}$$

for all sufficiently small s and all sufficiently large t . To show this, notice that for all $1 \leq k \leq (n - 1)$, and all $s > 0$,

$$\left| \frac{\partial^k}{\partial s^k} f(s, t) + \Lambda t \tilde{F}^{(k)}(s) \right| \leq \Lambda \mathbb{E}[\tau^{k+1}] + \Lambda t \mathbb{E}[\tau^k \mathbf{1}(\tau \geq t)]. \tag{B.17}$$

For $k = n$

$$\left| \frac{\partial^n}{\partial s^n} f(s, t) + \Lambda t \tilde{F}^{(n)}(s) \right| \leq \Lambda \mathbb{E}[\tau^{n+1} \mathbf{1}(\tau < t)] + \Lambda t \mathbb{E}[\tau^n \mathbf{1}(\tau \geq t)] = o(t), \tag{B.18}$$

as $t \rightarrow \infty$. Hence, since $|\tilde{F}^{(k)}(s)| \leq \mathbb{E}\tau^k < \infty$, $s \geq 0$, $0 \leq k \leq n$, (B.17) and (B.18) imply that

$$\left| \frac{\partial^k}{\partial s^k} f(s, t) \right| = O(t), \tag{B.19}$$

uniformly for all $1 \leq k \leq n$, and all $s > 0$. Furthermore, for $k = 1$, for any $\epsilon > 0$, $\epsilon < \rho$,

$$(1 - \epsilon) \rho t \leq \frac{\partial}{\partial s} f(s, t) \leq (1 + \epsilon) \rho t, \tag{B.20}$$

for all sufficiently small s and all sufficiently large t . Using (B.19), after some straightforward algebra, which we skip here, we arrive at

$$\left| (-1)^k e^{f(s,t)} \frac{\partial^k}{\partial s^k} e^{-f(s,t)} - \left(\frac{\partial}{\partial s} f(s, t) \right)^k \right| = O(t^{k-1}), \tag{B.21}$$

uniformly in $s > 0$. Subsequently, combination of (B.20) and (B.21) yields (B.16) (with possibly two different ϵ in (B.20) and (B.16)). Finally, (B.15) and (B.16), give (B.6), for $1 \leq k \leq n$. This finishes the proof of Lemma B.1.

Appendix C. Finite number of subexponential on-off processes: upper bound

For a finite number of long-tailed on-off processes we can obtain the following general upper bound.

Theorem C.1. *Let $F_1 \in \mathcal{S}$. If $r > c/N$, then*

$$\lim_{x \rightarrow \infty} \frac{\mathbb{P}[Q^N > x]}{\int_{x/(r-c_N)}^{\infty} \mathbb{P}[\tau^{\text{on}} > u] du} \leq NK_N, \tag{C.1}$$

where $c_N = c/N$, and K_N is given by Theorem 4.3, Equation (4.10), with c_N in place of c .

Proof. As in the proof of Theorem 4.4, we have the decomposition

$$\begin{aligned} W_t &= \int_{-t}^0 (A_u^N - c) du = \sum_{i=1}^N \int_{-t}^0 (a_u^i - c/N) du \\ &\stackrel{\text{def}}{=} \sum_{i=1}^N w_t^i. \end{aligned}$$

Hence,

$$\begin{aligned} \mathbb{P}[Q_t^N > x] &= \mathbb{P}\left[\sup_{t \geq 0} \sum_{i=1}^N w_t^i > x\right] \\ &\leq \mathbb{P}\left[\sum_{i=1}^N \sup_{t \geq 0} w_t^i > x\right] \\ &\sim N \mathbb{P}\left[\sup_{t \geq 0} w_t^1 > x\right] \quad \text{as } x \rightarrow \infty \end{aligned} \tag{C.2}$$

$$= NK_N, \tag{C.3}$$

where (C.2) follows from Lemma A.4(i), since $\sup_{t \geq 0} w_t^i, 1 \leq i \leq N$, are i.i.d. subexponential (by Theorem 4.3) random variables; (C.3) follows also from Theorem 4.3. This proves (C.1).

Acknowledgements

The authors are very grateful to Onno Boxma for the preprint of his paper and to Claudia Klüppelberg for her helpful comments.

References

- [1] ABATE, J., CHOUDHURY, G. L. AND WHITT, W. (1994). Waiting-time tail probabilities in queues with long-tail service-time distributions. *Queueing Systems* **16**, 311–338.
- [2] ANANTHARAM, V. (1995). On the sojourn time of sessions at an ATM buffer with long-range dependent input traffic. *Proc. 34th IEEE CDC*. IEEE Control Systems Society, New York.
- [3] ANICK, D., MITRA, D. AND SONDHI, M. M. (1982). Stochastic theory of a data handling system with multiple sources. *Bell Syst. Techn. J.* **61**, 1871–1894.
- [4] ASMUSSEN, S. (1987). *Applied Probability and Queues*. Wiley, New York
- [5] ASMUSSEN, S. (1996). Rare events in the presence of heavy tails. In *Stochastic Networks: Stability and Rare Events*, eds. P. Glasserman, K. Sigman and D. D. Yao (Lecture Notes in Statist. **117**). Springer, New York.
- [6] ASMUSSEN, S., HENRIKSEN, L. F. AND KLÜPPELBERG, C. (1994). Large claims approximations for risk processes in a Markovian environment. *Stoch. Proc. Appl.* **54**, 29–43.
- [7] ASMUSSEN, S., SCHMIDL, H. AND SCHMIDT, V. (1997). Tail probabilities for non-standard risk and queueing processes with subexponential jumps. *Adv. Appl. Prob* **31**, 422–447.
- [8] ATHREYA, K. B. AND NEY, P. E. (1972). *Branching Processes*. Springer, New York.

- [9] BACCELLI, F. AND BREMAUD, P. (1994). *Elements of Queueing Theory: Palm-Martingale Calculus and Stochastic Recurrence*. Springer, New York.
- [10] BINGHAM, N. H. AND DONEY, R. A. (1974). Asymptotic properties of super-critical branching processes I: the Galton–Watson process. *Adv. Appl. Prob.* **6**, 711–731.
- [11] BINGHAM, N. H., GOLDIE, C. M. AND TEUGELS, J. L. (1987). *Regular Variation*. CUP, Cambridge.
- [12] BOXMA, O. J. (1996). Fluid queues and regular variation. *Perf. Eval.* **27&28**, 699–712.
- [13] BOXMA, O. J. (1997). Regular variation in a multi-source fluid queue. In *Teletraffic contributions for the Information Age*, Vol. 2 (Proc. ITC 15), ed. V. Ramaswami, P. E. Wirth. Elsevier, New York, pp. 391–402.
- [14] CHISTYAKOV, V. P. (1964). A theorem on sums of independent positive random variables and its application to branching random processes. *Theory Prob. Appl.* **9**, 640–648.
- [15] CHOUDHURY, G. L. AND WHITT, W. (1997). Long-tail buffer-content distributions in broadband networks. *Perf. Eval.* **30**, 177–190.
- [16] CINLAR, E. (1972). Superposition of point processes. In *Stochastic Point Processes: Statistical Analysis, Theory and Application*, ed. P. A. W. Lewis. Wiley, New York, pp. 594–606.
- [17] CINLAR, E. (1975). *Introduction to Stochastic Processes*. Prentice-Hall, Engleton Cliffs, NJ.
- [18] CLINE, D. B. H. (1986). Convolution tails, product tails and domains of attraction. *Prob. Theory Rel. Fields* **72**, 529–557.
- [19] CLINE, D. B. H. (1994). Intermediate regular and π variation. *Proc. London Math. Soc.* **68**, 594–616.
- [20] COHEN, J. W. (1973). Some results on regular variation for distributions in queueing and fluctuation theory. *J. Appl. Prob.* **10**, 343–353.
- [21] COHEN, J. W. (1974). Superimposed renewal processes and storage with gradual input. *Stoch. Proc. Appl.* **2**, 31–58.
- [22] COHEN, J. W. (1994). On the effective bandwidth in buffer design for the multi-server channels. Technical report, CWI Report BS-R9406, NL-1090 GB, Amsterdam, The Netherlands.
- [23] COX, D. R. AND SMITH, W. L. (1954). On the superposition of renewal processes. *Biometrika* **41**, 91–99.
- [24] DUFFIELD, N. G. AND O’CONNELL, N. (1995). Large deviations and overflow probabilities for the general single-server queue with applications. *Math. Proc. Camb. Phil. Soc.* **118**, 363–374.
- [25] ELWALID, A., HEYMAN, D., LAKSHMAN, T. V., MITRA, D. AND WEISS, A. (1995). Fundamental bounds and approximations for ATM multiplexers with applications to video teleconferencing. *IEEE J. Sel. Areas Commun.* **13**, 1004–1016.
- [26] EMBRECHTS, P., GOLDIE, C. M. AND VERAVERBEKE, N. (1979). Subexponentiality and infinite divisibility. *Z. Wahrscheinlichkeitsth.* **49**, 335–347.
- [27] GLYNN, P. V. AND WHITT, W. (1994). Logarithmic asymptotics for steady-state tail probabilities in a single-server queue. In *Studies in Applied Probability*, eds. J. Galambos and J. Gani (*J. Appl. Prob.* **31A**). Applied Probability Trust, Sheffield, England, pp. 131–156.
- [28] HEATH, D., RESNICK, S. AND SAMORODNITSKY, G. (1998). Heavy tails and long range dependence in on/off processes and associated fluid models. *Math. Operat. Res.* **23**, 145–165.
- [29] HEYMAN, D. P. AND LAKSHMAN, T. V. (1996). Source models for VBR broadcast-video traffic. *IEEE/ACM Trans. Networking* **4**, 40–48.
- [30] JELENKOVIĆ, P. R. AND LAZAR, A. A. (1996). Multiple time scale and subexponential asymptotic behavior of a network multiplexer. In *Stochastic Networks: Stability and Rare Events*, eds. P. Glasserman, K. Sigman and D. D. Yao. (Lecture Notes in Statist. **117**). Springer, New York, pp. 215–235.
- [31] JELENKOVIĆ, P. R. AND LAZAR, A. A. (1996). Multiplexing on-off sources with subexponential on periods. CTR Technical Report CU/CTR/TR 457-96-23, Columbia University, New York. ([www: http://www.ctr.columbia.edu/comet/publications.](http://www.ctr.columbia.edu/comet/publications))
- [32] JELENKOVIĆ, P. R. AND LAZAR, A. A. (1997). Multiplexing on-off sources with subexponential on periods: Part I. In *Proc. INFOCOM’97*. IEEE Computer Society, New York.
- [33] JELENKOVIĆ, P. R. AND LAZAR, A. A. (1997). Multiplexing on-off sources with subexponential on periods: Part II. In *Teletraffic contributions for the Information Age*, Vol. 2 (Proc. ITC 15), ed. V. Ramaswami, P. E. Wirth. Elsevier, New York, pp. 965–974.
- [34] JELENKOVIĆ, P. R. AND LAZAR, A. A. (1998). Subexponential asymptotics of a Markov-modulated random walk with queueing applications. *J. Appl. Prob.* **35**, 325–347.
- [35] JELENKOVIĆ, P. R., LAZAR, A. A. AND SEMRET, N. (1996). Multiple time scales and subexponentiality in MPEG video streams. In *Broadband communications: global infrastructure for the information age* (Proc. Internat. IFIP-IEEE Conference on Broadband Communications), ed. L. Mason and A. Casaca. Chapman and Hall, London.
- [36] JELENKOVIĆ, P. R., LAZAR, A. A. AND SEMRET, N. (1997). The effect of multiple time scales and subexponentiality of MPEG video streams on queueing behavior. *IEEE J. Sel. Areas Commun.* **15**, 1052–1071.
- [37] KARAMATA, J. (1930). Sur un mode de croissance régulière des fonctions. *Mathematica (Cluj)* **4**, 38–53.

- [38] KELLA, O. AND WHITT, W. (1992). A storage model with a two-state random environment. *Operat. Res.* **40**, 257–262.
- [39] KLÜPPELBERG, C. (1988). Subexponential distributions and integrated tails. *J. Appl. Prob.* **25**, 132–141.
- [40] KLÜPPELBERG, C. (1989). Subexponential distributions and characterizations of related classes. *Prob. Theory Rel. Fields* **82**, 259–269.
- [41] LELAND, W. E., TAQQU, M. S., WILLINGER, W. AND WILSON, D. V. (1993). On the self-similar nature of Ethernet traffic. In *SIGCOMM'93*, Assoc. for Computing Machinery, New York, pp. 183–193.
- [42] LIKHANOV, N., TSYBAKOV, B. AND GEORGANAS, N. D. (1995). Analysis of an ATM buffer with self-similar ('fractal') input traffic. In *INFOCOM'95*, IEEE Computer Society, New York, pp. 985–991.
- [43] LOYNES, R. M. (1968). The stability of a queue with non-independent inter-arrival and service times. *Proc. Camb. Phil. Soc.* **58**, 497–520.
- [44] NORROS, I. (1994). A storage model with self-similar input. *Queueing Systems* **16**, 387–396.
- [45] PAKES, A. G. (1975). On the tails of waiting-time distribution. *J. Appl. Prob.* **12**, 555–564.
- [46] PARULEKAR, M. AND MAKOWSKI, A. M. (1996). Tail probabilities for a multiplexer with self-similar traffic. In *INFOCOM'96*, IEEE Computer Society, New York.
- [47] PARULEKAR, M. AND MAKOWSKI, A. M. (1997). Tail probabilities for $M/G/\infty$ input processes (I): preliminary asymptotics. *Queueing Systems* **27**, 271–296.
- [48] RESNICK, S. AND SAMORODNITSKY, G. (1997). Performance decay in a single server queueing model with long range dependence. *Operat. Res.* **45**, 235–243.
- [49] ROLSKI, T., SCHLEGEL, S. AND SCHMIDT, V. (1999). Asymptotics of palm-stationary buffer content distribution in fluid flow queues. *Adv. Appl. Prob.* **31**, 235–253.
- [50] RUBINOVITCH, M. (1973). The output of a buffered data communication system. *Stoch. Proc. Appl.* **1**, 375–380.
- [51] RYU, B. K. AND LOWEN, S. B. (1996). Point process approaches to the modeling and analysis of self-similar traffic - part I: Model construction. In *INFOCOM'96*, IEEE Computer Society, New York.
- [52] WEISS, A. AND SHWARTZ, A. (1995). *Large Deviations for Performance Analysis: Queues, Communications, and Computing*. Chapman & Hall, New York.
- [53] WILLEKENS, E. AND TEUGELS, J. L. (1992). Asymptotic expansion for waiting time probabilities in an $M/G/1$ queue with long-tailed service time. *Queueing Systems* **10**, 295–312.