

Modeling and analysis of power-tail distributions via classical teletraffic methods

David Starobinski ^{a,*} and Moshe Sidi ^{b,**}

^a *Department of Electrical Engineering and Computer Sciences, University of California at Berkeley, Berkeley, CA 94720, USA*

E-mail: staro@eecs.berkeley.edu

^b *Department of Electrical Engineering, Technion – Israel Institute of Technology, Haifa 32000, Israel*

E-mail: moshe@ee.technion.ac.il

Received 25 January 1999; revised 24 February 2000

We propose a new methodology for modeling and analyzing power-tail distributions, such as the Pareto distribution, in communication networks. The basis of our approach is a fitting algorithm which approximates a power-tail distribution by a hyperexponential distribution. This algorithm possesses several key properties. First, the approximation can be achieved within any desired degree of accuracy. Second, the fitted hyperexponential distribution depends only on a few parameters. Third, only a small number of exponentials are required in order to obtain an accurate approximation over many time scales. Once equipped with a fitted hyperexponential distribution, we have an integrated framework for analyzing queueing systems with power-tail distributions. We consider the $GI/G/1$ queue with Pareto distributed service time and show how our approach allows to derive both quantitative numerical results and asymptotic closed-form results. This derivation shows that classical teletraffic methods can be employed for the analysis of power-tail distributions.

Keywords: communication networks, $GI/G/1$ queue, multiple time-scale traffic, fitting, heavy-tail distribution, hyperexponential distribution, mixture of exponentials

1. Introduction

Recent studies have revealed that network traffic exhibits burstiness over multiple time scales [15,22]. In many circumstances, *power-tail* probability distributions have been found appropriate for capturing this salient feature (see [19 and references therein]). A random variable X has a power-tail distribution if its complementary cumulative distribution function (ccdf) $\overline{F}(t)$ satisfies

$$\overline{F}(t) = \Pr\{X > t\} \sim ct^{-\alpha} \quad \text{as } t \rightarrow \infty, \quad (1)$$

* The research was performed while this author was a Ph.D. student at the Technion – Israel Institute of Technology.

** This work was supported by the Israel Science Foundation administrated by the Academy of Science and Humanities.

where α and c are positive constants, and $f(t) \sim g(t)$ means $\lim_{t \rightarrow \infty} f(t)/g(t) = 1$. The case $0 < \alpha \leq 1$ is usually not of practical interest in queueing analysis since $E\{X\} = \infty$ (in this work, we assume that $\alpha > 1$). The most encountered situation is $1 < \alpha \leq 2$ for which the random variable X has a finite mean but an infinite variance. Occurrence of such a distribution in the activity and/or silence period of an On/Off process gives rise to long-range dependence, i.e., a non-summable autocorrelation function [6]. A well-known power-tail distribution is the (translated) Pareto distribution for which

$$\bar{F}(t) = \Pr\{X > t\} = \frac{1}{(1 + at)^\alpha} \quad \text{for } t \geq 0 \text{ and } a > 0. \quad (2)$$

The Pareto distribution provides parsimonious modeling since it depends on only two parameters.

Unfortunately, power-tail distributions do not lend themselves to easy queueing analysis since their Laplace transforms are not explicit, in most cases (for an exception, see [5]). This explains why, so far, most of the queueing results involving power-tail distributions have only been obtained in asymptotic regimes (see [11 and references therein]). These asymptotic results have the merit of providing some insight into the relation between the power-tail distributions parameters and the queueing performance measures.

In order to obtain more quantitative results, several contributions have recently suggested to fit hyperexponential distributions, i.e., mixture of exponentials, to power-tail distributions [9,12,20] (see also the related works [2,21]). However, none of the fitting algorithms developed in these work provide a systematic way for deriving an approximation arbitrarily close to the original distribution. Moreover, the queueing results obtained via these approaches are only numerical.

Inspired by a work of Mandelbrot [17], we propose, here, a new methodology for fitting hyperexponential distributions to power-tail distributions. This new approach exhibits several advantages. First, the approximation can be made *arbitrarily close* to the exact distribution and bounds on the approximation error are easily obtained. Second, the fitted hyperexponential distribution depends only on a few parameters which are explicitly related to the parameters of the power-tail distribution. Third, only a small number of exponentials are required in order to obtain an accurate approximation over multiple time scales, e.g., a dozen of exponentials for five time scales. Once equipped with a fitted hyperexponential distribution, we have an integrated framework for analyzing queueing systems with power-tail distributions. We consider the $GI/G/1$ queue with Pareto distributed service time and show how our approach allows to derive both quantitative numerical results and asymptotic closed-form results.

This paper is organized as follows. In the next section, we present our fitting algorithm. We provide bounds on the approximation error and prove that the fitted hyperexponential distribution can be made arbitrarily close to the original power-tail distribution. As an illustration of the method, we provide an explicit expression for a hyperexponential distribution, termed pseudo-Pareto distribution, which can approxi-

mate arbitrarily closely the Pareto distribution. In section 3, we study the distribution of the waiting time in a queue with i.i.d. and arbitrarily distributed interarrival times and pseudo-Pareto distributed service time. We show that it is straightforward to obtain a numerical solution for the waiting time distribution, even when the number of exponentials is very large. Moreover, as the number of exponentials tends to infinity, we derive an analytical expression for the tail of the waiting time distribution. The last section is devoted to a summary of the work and concluding remarks.

2. The fitting algorithm

Our algorithm proceeds in two stages. The first and most significant stage focuses on fitting a mixture of exponentials to the behavior of the tail of the power-tail distribution. The second stage provides a fitting for small values of t and ensures that the mixture of exponentials is indeed a probability distribution. As an example, we consider the case of the Pareto distribution defined in equation (2).

2.1. Mimicking the long term behavior

Consider the function $R(t) = ct^{-\alpha}$. We want to derive an expression for a mixture of exponentials which can capture the behavior of $R(t)$ from some value of t and over an arbitrary large number of time scales. Our starting point relies on the fact that $ct^{-\alpha}$ is the Laplace transform of the function $r(s) = cs^{\alpha-1}/\Gamma(\alpha)$, where $\Gamma(\cdot)$ is the Gamma function. We can, therefore, express $R(t)$ in the following way:

$$R(t) = c \int_0^\infty \frac{s^{\alpha-1} e^{-st}}{\Gamma(\alpha)} ds. \tag{3}$$

In the sequel of this subsection, we let $c = 1$ since it is merely a constant of proportionality. The integral appearing in the right-hand side of equation (3) can be approximated by a Riemann sum. However, we know from the Tauberian theorems (see [10, pp. 442–448]) that the behavior of $R(t)$ for large values of t is closely related to the behavior of $r(s)$ near $s = 0$. The choice of a fixed grid would not be wise. It would put too much emphasis on large values of s which correspond to “high frequencies” and not enough on small values of s corresponding to “low frequencies”. We perform therefore the following change of variables from s to u , $s = B^{-u}$, where $B > 1$ is a parameter which controls the accuracy of the approximation, as is made clear later. We note that choosing a fixed grid for the variable u is equivalent to choosing a logarithmic grid for s . After the change of variables, equation (3) can be rewritten as

$$\begin{aligned} R(t) &= \frac{\log B}{\Gamma(\alpha)} \int_{-\infty}^\infty B^{-\alpha u} \exp(-tB^{-u}) du \\ &= \frac{\log B}{\Gamma(\alpha)} \sum_{n=-\infty}^\infty \int_{n-1/2}^{n+1/2} B^{-\alpha u} \exp(-tB^{-u}) du. \end{aligned} \tag{4}$$

Equation (4) can be approximated by a Riemann sum if we replace each integrand with its mid-span value. It turns out, however, that a better approximation can be obtained if only the exponent portion of the integrand is replaced with its mid-span value. We have then

$$\begin{aligned} R(t) &\approx \frac{\log B}{\Gamma(\alpha)} \sum_{n=-\infty}^{\infty} \exp(-tB^{-n}) \int_{n-1/2}^{n+1/2} B^{-\alpha u} du \\ &= \frac{B^{\alpha/2} - B^{-\alpha/2}}{\Gamma(\alpha + 1)} \sum_{n=-\infty}^{\infty} B^{-\alpha n} \exp(-tB^{-n}) \equiv R_1(t). \end{aligned} \quad (5)$$

As proven in the next subsection, with $B \rightarrow 1$, the approximation $R_1(t)$ can be made arbitrarily close to $R(t)$. The last step of the algorithm is to truncate the infinite sum $R_1(t)$ and approximate it by a finite sum $R_2(t)$, where

$$R_2(t) = \frac{B^{\alpha/2} - B^{-\alpha/2}}{\Gamma(\alpha + 1)} \sum_{n=M}^N B^{-\alpha n} \exp(-tB^{-n}). \quad (6)$$

The idea behind this truncation is the following. On one hand, values of n below M correspond to high frequencies which have almost no effect on the long-term behavior of $R(t)$. On the other hand, values of n larger than N correspond to very low frequencies (or very large values of t) falling beyond the scope of interest. Note that the approximation $R_2(t)$ is very parsimonious since it depends on only four parameters: α , B , M and N . As an illustration of the fitting method, we consider the example of a power-tail function $R(t) = t^{-3/2}$ with its approximating function $R_2(t)$. The values chosen for the parameters of $R_2(t)$ are $\alpha = 3/2$, $B = 2$, $M = 0$ and $N = 20$. As one can see from figure 1(a), the quality of the approximation is excellent over the whole domain $t \in [10, 10^5]$. The fitting is less tight for values of t outside this domain due to the reasons mentioned above.

2.2. Approximation errors and bounds

The objective of this subsection is to develop a procedure for bounding the approximation error. Based on the results of this procedure, we prove that $R_2(t)$ can approximate arbitrarily closely the exact function $R(t)$ over any interval $[t_a, t_b]$ ($0 < t_a < t_b < \infty$). Moreover, the bounds provide a very useful insight into the problem of setting the values of the parameters of $R_2(t)$.

We define the relative approximation error between $R(t)$ and $R_2(t)$ as

$$\text{Err}[R(t), R_2(t)] = \frac{|R(t) - R_2(t)|}{\min(R(t), R_2(t))}. \quad (7)$$

The procedure for bounding the error is based on the derivation of the following two functions: a function $R_{\text{up}2}(t)$ which bounds $R(t)$ and $R_2(t)$ from above and another

Figure 1(a): Original $R(t)$ and approximation $R_2(t)$

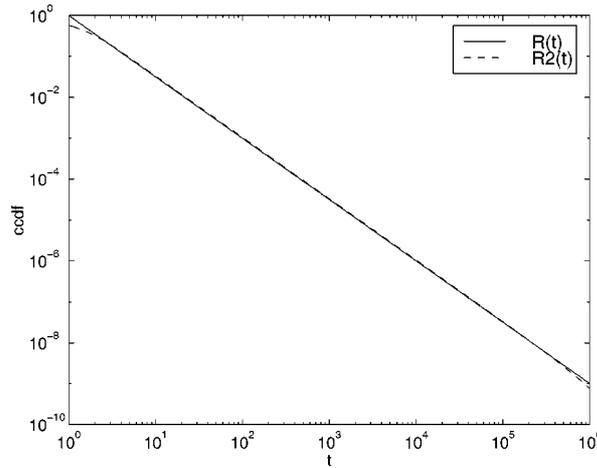


Figure 1(b): Approximation error and a bound

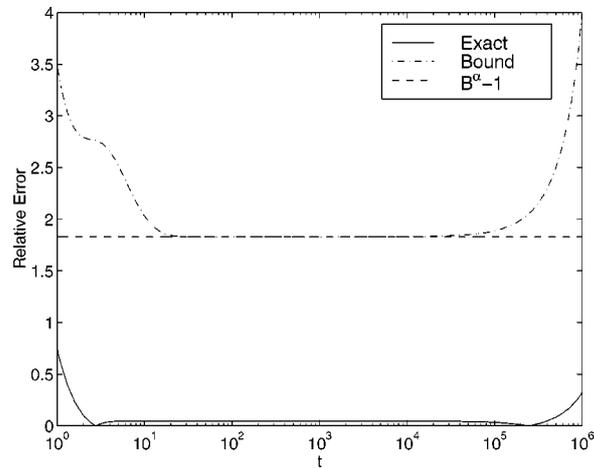


Figure 1. Example of fitting a mixture of exponentials to a power-tail function.

function $R_{lo2}(t)$ which bounds $R(t)$ and $R_2(t)$ from below. Once we obtain an expression for these functions, we readily get the following bound on the error:

$$\text{Err}[R(t), R_2(t)] \leq \frac{R_{up2}(t) - R_{lo2}(t)}{R_{lo2}(t)}. \tag{8}$$

The fitting method described in the previous subsection is based on two approximations: (i) the discretization of an integral and (ii) the truncation of an infinite sum. We begin by considering the error due to the discretization. We recall that the discretization leading to $R_1(t)$ has been obtained by replacing the exponent portion $f(u) = \exp(-tB^{-u})$ of each integrand in equation (4) with its mid-span value $f(n) = \exp(-tB^{-n})$. The function $f(u)$ is strictly increasing with u . Therefore, if

we replace $f(u)$ with $f(n + 1/2)$, which is the value of the function at the right-most point of the integration interval, we obtain an upper bound $R_{\text{up}1}(t)$ on both $R(t)$ and $R_1(t)$. The expression for $R_{\text{up}1}(t)$ is

$$R_{\text{up}1}(t) = \frac{B^{\alpha/2} - B^{-\alpha/2}}{\Gamma(\alpha + 1)} \sum_{n=-\infty}^{\infty} B^{-\alpha n} \exp(-tB^{-n-1/2}). \quad (9)$$

In a strictly analogous way, we obtain a lower bound $R_{\text{lo}1}(t)$ on both $R(t)$ and $R_1(t)$ by replacing $f(u)$ with $f(n - 1/2)$:

$$R_{\text{lo}1}(t) = \frac{B^{\alpha/2} - B^{-\alpha/2}}{\Gamma(\alpha + 1)} \sum_{n=-\infty}^{\infty} B^{-\alpha n} \exp(-tB^{-n+1/2}). \quad (10)$$

Before we proceed, it is instructive to derive a bound on the approximation error between $R(t)$ and $R_1(t)$:

$$\begin{aligned} \text{Err}[R(t), R_1(t)] &\leq \frac{R_{\text{up}1}(t)}{R_{\text{lo}1}(t)} - 1 = \frac{\sum_{n=-\infty}^{\infty} B^{-\alpha n} \exp(-tB^{-n-1/2})}{\sum_{n=-\infty}^{\infty} B^{-\alpha n} \exp(-tB^{-n+1/2})} - 1 \\ &= B^\alpha \frac{\sum_{\tilde{n}=-\infty}^{\infty} B^{-\alpha \tilde{n}} \exp(-tB^{-\tilde{n}+1/2})}{\sum_{n=-\infty}^{\infty} B^{-\alpha n} \exp(-tB^{-n+1/2})} - 1 = B^\alpha - 1. \end{aligned} \quad (11)$$

From equation (11), we can already draw three intermediate conclusions. First, the bound on the approximation error between $R(t)$ and $R_1(t)$ is independent of t . Second, the approximation error can be made arbitrarily small by letting B approach 1 from above. Finally, we observe that as α becomes larger, smaller values of B will be required for achieving the same degree of accuracy.

The next stage consists of bounding the approximation error due to the discretization and the truncation, altogether. Let us first consider the derivation of the lower bound. We define $R_{\text{lo}2}(t)$ as the truncation of $R_{\text{lo}1}(t)$

$$R_{\text{lo}2}(t) = \frac{B^{\alpha/2} - B^{-\alpha/2}}{\Gamma(\alpha + 1)} \sum_{n=M}^N B^{-\alpha n} \exp(-tB^{-n+1/2}). \quad (12)$$

Using the same arguments as in the previous paragraph, it follows that $R_{\text{lo}2}(t)$ is a lower bound on $R_2(t)$. At the same time, $R_{\text{lo}2}(t)$ is also a lower bound on $R(t)$ since $R_{\text{lo}2}(t) < R_{\text{lo}1}(t) < R(t)$. The derivation of the upper bound $R_{\text{up}2}(t)$ is more lengthy. For this reason, we defer the technical details to appendix A, and provide directly the final expression for this function:

$$\begin{aligned} R_{\text{up}2}(t) &= \frac{B^{\alpha/2} - B^{-\alpha/2}}{\Gamma(\alpha + 1)} \sum_{n=M}^N B^{-\alpha n} \exp(-tB^{-n-1/2}) \\ &\quad + \frac{\exp(-tB^{1/2-M})}{\Gamma(\alpha)} \sum_{k=0}^{\lceil \alpha-1 \rceil} \frac{\lceil \alpha-1 \rceil!}{k!} \frac{B^{-k(1/2-M)}}{t^{\lceil \alpha-1 \rceil - k + 1}} + \frac{B^{-\alpha(N+1/2)}}{\Gamma(\alpha + 1)}. \end{aligned} \quad (13)$$

The first term in equation (13) refers to the discretization error. The second term provides an upper bound on the truncation error due to the cut-off of high frequencies. Note that this term has no influence on the long-term behavior, since it decays exponentially fast. Moreover, for any $t > 0$, it can be made as small as desired by letting $M \rightarrow -\infty$. The third term is an upper bound on the low frequency error. This term vanishes if one lets $N \rightarrow \infty$. We conclude that $R_2(t)$ can approximate arbitrarily closely the exact function $R(t)$, over any finite interval, by letting B approach 1, and taking M small enough and N large enough. For typical values of α , i.e., $1 < \alpha < 2$, our experience teaches that values of B , M and N ranging, respectively, from 2 to 3, -1 to 1, and 15 to 25, yield a very good approximation over four or five time scales.

To illustrate the results of this procedure, let us consider again the example given at the end of the last subsection ($R(t) = t^{-3/2}$ and $R_2(t)$ with parameters $\alpha = 3/2$, $B = 2$, $M = 0$ and $N = 20$). In figure 1(b), we plot the approximation error and a bound on it. The exact expression for the approximation error is given by equation (7). The expression for the bound is obtained by substituting (12) and (13) into the right-hand side of equation (8). We observe that the bound is nearly constant in the mid-frequency region and approximately equals to $B^\alpha - 1$. This is expected since the discretization error is the main source of inaccuracy in this region (see equation (11)). We remark also that the bound is not very close to the actual error which is very small within the domain of interest. Nevertheless, the *qualitative* behavior is similar. This similarity provides useful guidelines for setting the values of the parameters of $R_2(t)$. For instance, let us assume that we want to obtain a “good” estimate of $R(t)$ over the interval $t \in [10, 10^5]$ (by *good*, we mean an approximation which is not influenced by the truncation). Computations show that the bound on the high-frequency error, given by the second term in the right-hand side of equation (13), is smaller than $R(t)$ by at least one order of magnitude when $t > 10$. The bound on the low-frequency error, given by the third term in the right-hand side of equation (13), becomes significant with regard to $R(t)$ only when $t > 10^5$. The selected values for M and N are therefore reasonable. Finally, note that the approximation can be improved by letting B approach 1. However, this will require in turn to decrease the value of M and increase the value of N .

2.3. Matching a probability distribution

Our goal in this subsection is to show how a hyperexponential distribution can be fitted to a power-tail probability distribution. For this purpose, we propose to match the Pareto distribution defined in equation (2). Using the Laplace transform representation, the ccdf of the Pareto distribution can be expressed as

$$\bar{F}(t) = \frac{1}{(1+at)^\alpha} = \int_0^\infty \frac{s^{\alpha-1} e^{-s(at+1)}}{\Gamma(\alpha)} ds. \quad (14)$$

A probabilistic interpretation of equation (14) is that the Pareto distribution represents a mixture of exponential distributions where the parameter of the exponential distribution is gamma distributed (see [13, p. 233]).

Following the same steps as in section 2.1, we obtain the following approximation for $\bar{F}(t)$:

$$\bar{F}(t) \approx \frac{B^{\alpha/2} - B^{-\alpha/2}}{\Gamma(\alpha + 1)} \sum_{n=M}^N B^{-\alpha n} \exp(-B^{-n}) \exp(-aB^{-n}t). \quad (15)$$

The expression in the right-hand side of equation (15) must be slightly modified in order to obtain a probability distribution. We define the sum of the coefficients of the exponentials in that expression as

Figure 2(a): Complementary cumulative distribution

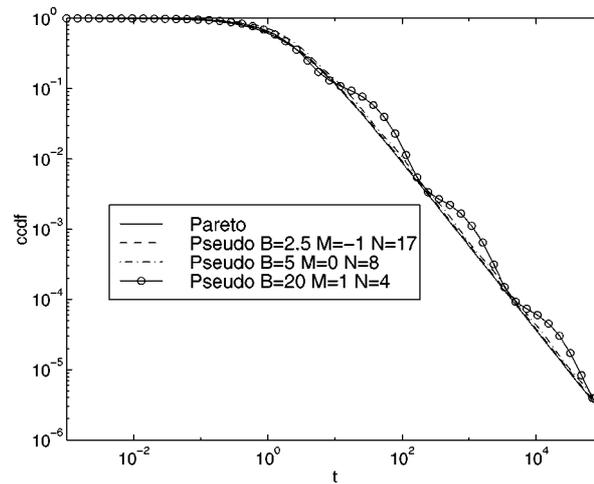


Figure 2(b): Relative approximation error

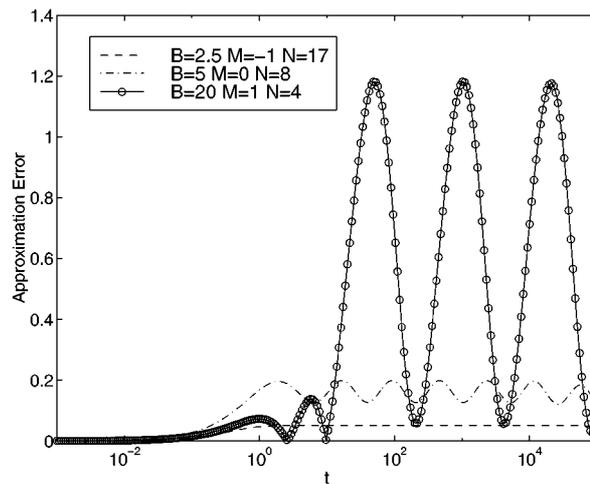


Figure 2. Comparison between a Pareto ccdf and three fitted pseudo-Pareto ccdfs consisting, respectively, of five, ten, and twenty exponentials.

$$\omega = \frac{B^{\alpha/2} - B^{-\alpha/2}}{\Gamma(\alpha + 1)} \sum_{n=M}^N B^{-\alpha n} \exp(-B^{-n}). \quad (16)$$

The following expression corresponds to a hyperexponential distribution:

$$\begin{aligned} \bar{G}(t) &= (1 - \omega) \exp(-aB^{-(M-1)}t) \\ &+ \frac{B^{\alpha/2} - B^{-\alpha/2}}{\Gamma(\alpha + 1)} \sum_{n=M}^N B^{-\alpha n} \exp(-B^{-n}) \exp(-aB^{-n}t). \end{aligned} \quad (17)$$

For large values of B , it may happen that ω is larger than 1 (due to the discretization error). In such a case, the value of M must be appropriately increased in order to ensure that ω is smaller than 1. Note that an alternative way for deriving an hyperexponential distribution is simply to divide (15) by (16). A potential drawback of this approach is that the inaccuracy resulting from the truncation of high frequencies may have an impact on the quality of the approximation for large values of t .

Using the same terminology as [20], we refer to $G(t)$ as a *pseudo-Pareto* distribution. Of course, the pseudo-Pareto distribution can be made arbitrarily close to the exact Pareto distribution by letting $B \rightarrow 1$, $N \rightarrow \infty$ and $M \rightarrow -\infty$. As an illustration of the fitting method, a Pareto distribution with cdf $\bar{F}(t) = 1/(1 + 0.5 \cdot t)^{1.2}$ is compared, in figure 2 to three fitted pseudo-Pareto distributions consisting, respectively, of five, ten, and twenty exponentials. From figure 2(a), we see that as few as five exponentials are enough in order to mimic the behavior of the Pareto distribution over several time scales. Of course, increasing the number of exponentials leads to a smaller approximation error, as one can observe from figure 2(b). For instance, using a pseudo-Pareto distribution consisting of twenty exponentials yields an approximation error smaller than 5% within the domain of interest.

3. GI/G/1 queueing analysis

In this section, we consider the analysis of the (actual) waiting time W in a GI/G/1 queue with Pareto distributed service time (note that W is closely related to the buffer content of a queue fed by an On/Off fluid process with Pareto distributed activity period and arbitrarily distributed silence period [4]). In order to analyze this queueing system, it is necessary to have the Laplace transform of the Pareto distribution. Since there is evidently no convenient expression for the Laplace transform of the Pareto distribution, our approach is to model the Pareto distribution with a pseudo-Pareto distribution. Once equipped with a pseudo-Pareto distribution, we show, in the sequel, that it is straightforward to derive a numerical solution for the waiting time distribution, even when the number of exponentials is very large. Moreover, we show that an asymptotic closed-form expression for $\Pr(W > t)$ prevails as $N \rightarrow \infty$ and $t \rightarrow \infty$. This expression is shown to coincide with a well-known result of Pakes [18] (see also [7]), as $B \rightarrow 1$.

Using the notion of weak convergence [3], Feldmann and Whitt [9] proved that it is theoretically possible to approximate arbitrarily closely the waiting time distribution in a $GI/G/1$ queue with a Pareto service time distribution by the waiting distribution in a $GI/G/1$ queue with a hyperexponential service time distribution. One of our main contributions, here, is to show how such an arbitrarily close approximation can be achieved in practice, via the pseudo-Pareto distribution.

3.1. Theoretical results

We consider a $GI/G/1$ queue with arrival rate λ and pseudo-Pareto service time distribution with mean $1/\mu$. The service policy of the queue is FIFO. We define the load as $\rho = \lambda/\mu$ and assume that it is smaller than 1. We denote the Laplace transforms of the interarrival and service time distributions, respectively, by $A^*(s)$ and $G^*(s)$. In the case of the pseudo-Pareto distribution, the expression for $G^*(s)$ is

$$G^*(s) = \frac{(1-\omega)aB^{-(M-1)}}{s+aB^{-(M-1)}} + \frac{B^{\alpha/2} - B^{-\alpha/2}}{\Gamma(\alpha+1)} \sum_{n=M}^N \frac{aB^{-(\alpha+1)n} \exp(-B^{-n})}{s+aB^{-n}}. \quad (18)$$

We observe that $G^*(s)$ is a rational function, with denominator of degree $N - M + 2$. The class of $GI/G/1$ queues with service time distribution having a rational Laplace transform is studied in [8, pp. 322–329]. The Laplace transform of the waiting time probability distribution is given by the following formula (see [8, equation (5.190)]):

$$W^*(s) = \prod_{n=M-1}^N \frac{-\sigma_n(s+aB^{-n})}{aB^{-n}(s-\sigma_n)}, \quad (19)$$

where σ_n , $n \in \{M-1, M, \dots, N-1, N\}$, correspond to the roots (zeroes) in the left half-plane $\text{Re}(s) < 0$ of the function $\Delta^*(s) = -1 + A^*(-s)G^*(s)$. The main computational effort required in order to invert $W^*(s)$ is the determination of the roots of $\Delta^*(s)$. The following proposition reduces considerably this effort.

Proposition 1. If $\rho < 1$, then $\Delta^*(s)$ has $N - M + 2$ distinct real roots in the left half-plane $\text{Re}(s) < 0$. A unique root, denoted by σ_n , is contained in each interval $(-aB^{-n}, -aB^{-(n+1)})$, where $n \in \{M-1, M, \dots, N-2, N-1\}$. An additional root σ_N is located in the interval $(-aB^{-N}, 0)$.

Proof. See appendix B. □

Since the roots are known to be real and to belong to distinct intervals, it is very easy to determine them with any simple search procedure. Note, also, that the scope of the above theorem can be easily extended to general hyperexponential distributions.

As $N \rightarrow \infty$, explicit asymptotic results on the location of σ_n can be obtained for large (positive) values of n . In such a case, it turns out that the location of σ_n gets very close to $s = -aB^{-n}$. We guess, therefore, that $\sigma_n = -aB^{-n} + \gamma_n$, where

$\gamma_n = o(B^{-n})$ (the notation $f(n) = o(g(n))$ means $\lim_{n \rightarrow \infty} f(n)/g(n) = 0$). We substitute this guess into $\Delta^*(s)$ and solve for γ_n . We obtain the following result:

Proposition 2. Let $N \rightarrow \infty$ and

$$\xi_n = -aB^{-n} + \frac{\lambda(B^{\alpha/2} - B^{-\alpha/2})}{(1 - \rho)\Gamma(\alpha + 1)} B^{-\alpha n}.$$

Then,

- (i) $\Delta^*(\xi_n) = o(B^{-n})$;
- (ii) $\sigma_n = \xi_n + o(B^{-\alpha n})$.

Proof. See appendix C. □

The significance of the above proposition is three-fold. First, it provides a good starting point for the root search procedure. Second, for large values of n , it states that ξ_n represents a very accurate approximation for σ_n . Third, it allows to obtain an asymptotic result on the waiting time distribution, as shown next.

In order to obtain an asymptotic expression for $\Pr(W > t)$, we perform a partial fraction expansion of $W^*(s)$,

$$W^*(s) = \prod_{n=M-1}^N \frac{-\sigma_n(s + aB^{-n})}{aB^{-n}(s - \sigma_n)} = \prod_{n=M-1}^N \frac{-\sigma_n}{aB^{-n}} + \sum_{n=M-1}^N \frac{\nu_n}{s - \sigma_n}. \quad (20)$$

The values of the coefficients ν_n are easily computed by resorting to the residue theorem. Moreover, as $N \rightarrow \infty$, it can be shown that

$$\nu_n = aB^{-n} + \xi_n + o(B^{-\alpha n}) = \frac{\lambda(B^{\alpha/2} - B^{-\alpha/2})}{(1 - \rho)\Gamma(\alpha + 1)} B^{-\alpha n} + o(B^{-\alpha n}). \quad (21)$$

This result is proven in appendix D using the asymptotic expression for σ_n given by proposition 2. The general expression for $\Pr(W > t)$ is then

$$\Pr(W > t) = \sum_{n=M-1}^N -\frac{\nu_n}{\sigma_n} e^{\sigma_n t} \quad \text{for } t > 0. \quad (22)$$

As we saw in section 2, the long-term behavior of $\Pr(W > t)$ is determined by the elements of the sum with large index n . The knowledge of the asymptotic behavior of ν_n and σ_n allows to obtain an analytical expression for $\Pr(W > t)$ as $N \rightarrow \infty$ and $t \rightarrow \infty$.

Proposition 3. As $N \rightarrow \infty$ and $t \rightarrow \infty$, one has

$$\Pr(W > t) \sim \frac{\lambda(B^{\alpha/2} - B^{-\alpha/2})}{a(1 - \rho)\Gamma(\alpha + 1)} \sum_{n=0}^N B^{(-\alpha+1)n} \exp(-aB^{-n}t). \quad (23)$$

Proof. See appendix E. □

Using the same kind of reasoning as in section 2, we know that, as $N \rightarrow \infty$, $t \rightarrow \infty$ and $B \rightarrow 1$,

$$\frac{B^{(\alpha-1)/2} - B^{(-\alpha+1)/2}}{\Gamma(\alpha)} \sum_{n=0}^N B^{(-\alpha+1)n} \exp(-aB^{-n}t) \sim (at)^{-\alpha+1}. \quad (24)$$

Combining equation (24) with equation (23), we obtain

$$\Pr(W > t) \sim \frac{\lambda a^{-\alpha}}{(1-\rho)(\alpha-1)} t^{-\alpha+1}, \quad (25)$$

as $N \rightarrow \infty$, $t \rightarrow \infty$ and $B \rightarrow 1$. This relation corresponds to the formula of Pakes [18] which states that the waiting time ccdf in a $GI/G/1$ queue satisfies equation (25) when the service time has a power-tail ccdf $\bar{F}(t) \sim (at)^{-\alpha}$ (actually, equation (25) is only a special case of Pakes' formula which applies also to more general subexponential distributions).

3.2. Numerical results

We present some numerical illustrations of the theoretical results derived in the previous subsection. We consider an $M/G/1$ queue with arrival rate $\lambda = 0.2$ and Pareto service time distribution with ccdf $\bar{F}(t) = 1/(1+2t)^{1.5}$ and mean $1/\mu = 1$. In figure 3, we compare numerical approximated results with simulated results for the waiting time distribution. The numerical results are obtained by replacing the Pareto distribution with a pseudo-Pareto distribution with parameters $B = 2$, $M = -1$ and $N = 25$. The derivation of the waiting time distribution is performed by first evaluating the quantities σ_n and ν_n , and then substituting their values into equation (22). The simulated results are obtained using the BONE's network simulator. Each simulation lasts 10^8 time units and 20 independent replications are run. The simulated results are presented with 99% confidence intervals. Figure 3 shows excellent agreement between the approximated and simulated results. Note also that, in this example, the asymptotic expressions for σ_n and ν_n (provided by proposition 2 and equation (21)) differ from the exact values for these quantities by less than 1% when $n \geq 15$.

In the next example, we compare the performance of $GI/G/1$ queues all having the same mean arrival rate and service time distribution but with different interarrival distributions. We assume that the service time distribution follows the pseudo-Pareto distribution described in the previous paragraph. Regarding the interarrival time, we consider three different distributions, all with mean $1/\lambda = 5$, that is a 2-stage Erlangian distribution, an exponential distribution and a hyperexponential distribution with ccdf $0.1 \exp(-0.05t) + 0.9 \exp(-0.3t)$. In each case, the derivation of the waiting time distribution is performed by first computing the quantities σ_n and ν_n , and then using equation (22). In figure 4, we present the behavior of the waiting time distribution in $GI/G/1$ queues with the mentioned interarrival time distributions. We observe that for small values of t , the waiting time distribution is strongly dependent on the

Validation of Numerical Results - 20 Replications

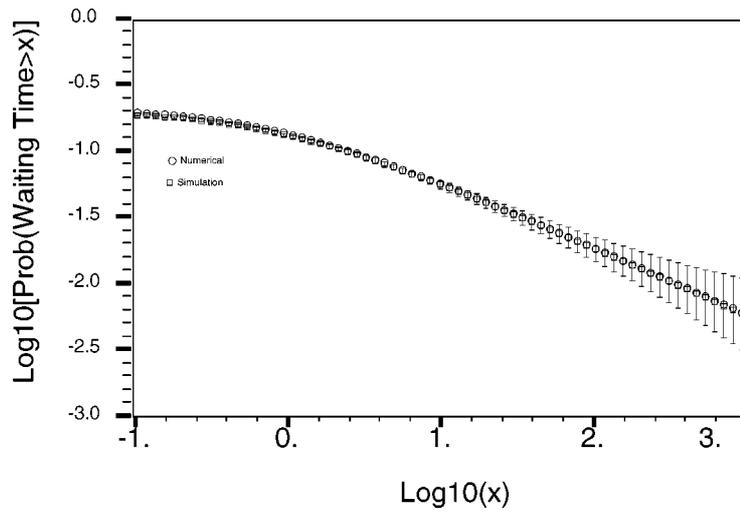


Figure 3. Waiting time distribution in an $M/G/1$ queue with Pareto service distribution: numerical approximated results versus simulated results.

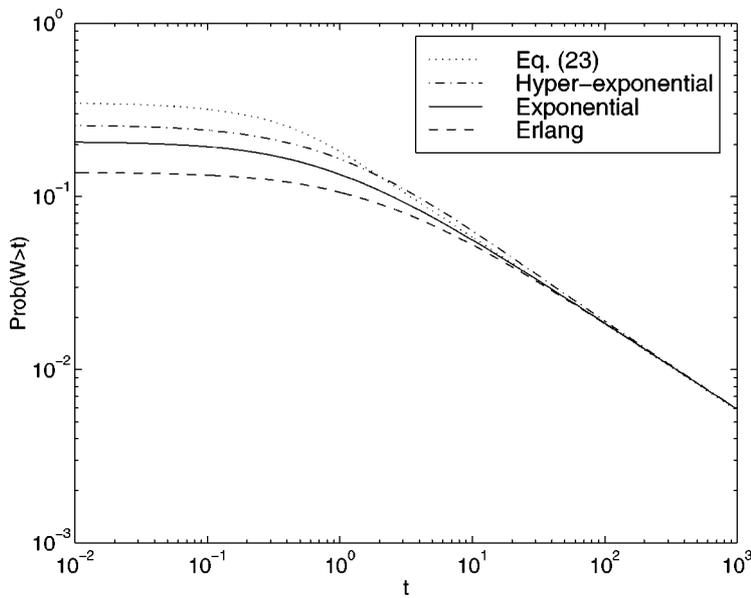


Figure 4. Waiting time distribution in $GI/G/1$ queues with pseudo-Pareto service distribution and various interarrival distributions, and comparison with asymptotic results.

interarrival time distribution. The probability of waiting, $\Pr(W > 0)$, is the largest for the hyperexponential distribution and the smallest for the Erlangian distribution. From equation (23), we also obtain an asymptotic expression for the waiting time distribution

(we take $N = 25$). This expression is also depicted in figure 4. Figure 4 indicates that the asymptotic region corresponds to rather large values of delay, i.e., $t > 100$. This result is in agreement with earlier work, see, e.g., [1], which already noticed that subexponential asymptotics may not provide especially good approximations.

4. Concluding remarks

In this work, we showed that classical teletraffic methods can be employed for the modeling and analysis of power-tail distributions in queueing systems. From the modeling point of view, we introduced a new algorithm which fits a hyperexponential distribution to a power-tail distribution. The fitted hyperexponential distribution depends only on a few parameters and provides parsimonious modeling. Also, the parameters of the power-tail distribution appear explicitly in the expression for the hyperexponential distribution. We showed that the approximation can be obtained within any desired degree of accuracy. As an example, we derived a new distribution, termed pseudo-Pareto distribution, which represents the “hyperexponential” counterpart of the Pareto distribution. From the analysis point of view, we considered the $GI/G/1$ queue and showed that when modeling the service time distribution with a pseudo-Pareto distribution, both quantitative numerical results and asymptotic closed-form results can be obtained. Our methodology provides new insight into the impact of system characteristics, such as the interarrival time distribution, on performance measures. Moreover, it enables to state the domain of validity of asymptotic results.

Our analysis of the $GI/G/1$ queue was based on the inversion of a Laplace transform. This approach has the advantage of resulting in an explicit expression for the waiting time distribution. Note that since the pseudo-Pareto distribution belongs to the family of phase-type distributions, matrix analytical methods [14,16] could also have been employed for deriving numerical results. In addition, matrix analytic methods can be used for analyzing the pseudo-Pareto distribution in a variety of queueing models.

We conclude with the following remarks. In section 3, we assumed that the service time distribution follows the pseudo-Pareto distribution. In fact, the results obtained in that section can be generalized to other hyperexponential distributions fitting power-tail distributions. Next, the bound on the approximation error derived in section 2.2 served to prove that a power-tail distribution can be approximated arbitrarily closely by a hyperexponential distribution. It provided also useful guidelines for the setting of the parameters M and N . Future work may look for a tighter bound on the approximation error which would give more quantitative insight into the setting of the parameter B . Finally, we note that the results of this work can be employed for obtaining upper bounds in network of queues, using the network calculus for “sums of exponentials” developed in [23].

Appendix A. Proof of equation (13)

In this appendix, we derive an expression for $R_{\text{up}2}(t)$. For this purpose, we start from the expression for $R(t)$ given by equation (5) and divide the integral into three

parts

$$\begin{aligned}
 R(t) &= \int_0^\infty \frac{s^{\alpha-1} e^{-st}}{\Gamma(\alpha)} ds \\
 &= \underbrace{\int_0^{B^{-1/2-N}} \frac{s^{\alpha-1} e^{-st}}{\Gamma(\alpha)} ds}_{(I)} + \underbrace{\int_{B^{-1/2-N}}^{B^{1/2-M}} \frac{s^{\alpha-1} e^{-st}}{\Gamma(\alpha)} ds}_{(II)} + \underbrace{\int_{B^{1/2-M}}^\infty \frac{s^{\alpha-1} e^{-st}}{\Gamma(\alpha)} ds}_{(III)}. \quad (26)
 \end{aligned}$$

The consequence of the truncation is to neglect parts (I) and (III) of equation (26). Part (I) in equation (26) represents the contribution of low frequencies, or correspondingly, values of u larger than $N + 1/2$ (remember that $s = B^{-u}$). A trivial upper bound on this integral can be derived as follows:

$$\int_0^{B^{-1/2-N}} \frac{s^{\alpha-1} e^{-st}}{\Gamma(\alpha)} ds < \int_0^{B^{-1/2-N}} \frac{s^{\alpha-1}}{\Gamma(\alpha)} ds = \frac{B^{-\alpha(N+1/2)}}{\Gamma(\alpha + 1)}. \quad (27)$$

From equation (27), we see that the contribution of low frequencies are negligible as long as $R(t) \gg B^{-\alpha N}$. Clearly the low-frequency error vanishes, as $N \rightarrow \infty$. Part (III) in equation (26) represents the contribution of high frequencies (values of u smaller than $M - 1/2$). An upper bound on this part can also be derived. For simplicity of exposition, we assume here that $M \leq 0$. We have then

$$\begin{aligned}
 \int_{B^{1/2-M}}^\infty \frac{s^{\alpha-1} e^{-st}}{\Gamma(\alpha)} ds &\leq \int_{B^{1/2-M}}^\infty \frac{s^{[\alpha-1]} e^{-st}}{\Gamma(\alpha)} ds \\
 &= \frac{\exp(-tB^{1/2-M})}{\Gamma(\alpha)} \sum_{k=0}^{[\alpha-1]} \frac{[\alpha-1]!}{k!} \frac{B^{-k(1/2-M)}}{t^{[\alpha-1]-k+1}}. \quad (28)
 \end{aligned}$$

The inequality follows from the fact that $s^{\alpha-1} \leq s^{[\alpha-1]}$ for $s \geq 1$ ($[x]$ denotes the smallest integer larger than or equal to x). From equation (28), we see that the high-frequency error has no influence on the long-term behavior since it decays exponentially fast with t (note that this property holds also when $M > 0$). Moreover, for any (fixed) value of t , the error can be made arbitrarily small by letting $B^{1/2-M} \rightarrow \infty$ or accordingly $M \rightarrow -\infty$. Finally, the discretization of part (II) in equation (26) leads to $R_2(t)$. An upper bound valid on both $R_2(t)$ and part (II) is easily obtained by resorting to the approach described in the previous paragraph. We have then

$$\begin{aligned}
 \int_{B^{-1/2-N}}^{B^{1/2-M}} \frac{s^{\alpha-1} e^{-st}}{\Gamma(\alpha)} ds &= \frac{\log B}{\Gamma(\alpha)} \sum_{n=M}^N \int_{n-1/2}^{n+1/2} B^{-\alpha u} \exp(-tB^{-u}) du \\
 &\leq \frac{B^{\alpha/2} - B^{-\alpha/2}}{\Gamma(\alpha + 1)} \sum_{n=M}^N B^{-\alpha n} \exp(-tB^{-n-1/2}). \quad (29)
 \end{aligned}$$

Summing (29), (28) and (27), we obtain the final expression for $R_{\text{up}2}(t)$ given by equation (13).

Appendix B. Proof of proposition 1

The fact that $\Delta^*(s)$ has $N - M + 2$ roots in the left-hand plane $\text{Re}(s) < 0$ is proven in [8, p. 323], by making use of Rouché’s theorem. In order to find the location of these roots, we study the behavior of the function $\Delta^*(s)$ for real, negative, values of s . We note that for such values of s , the function $A^*(-s)$ is continuous (since it is analytic), positive and bounded, i.e., $0 < A^*(-s) < 1$. The function $\Delta^*(s)$ has the same $N - M + 2$ points of discontinuities as $G^*(s)$. These points are located at $s = -aB^{-n}$, where $n \in \{M - 1, M, \dots, N\}$. In each interval $(-aB^{-n}, -aB^{-(n+1)})$, where $n \in \{M - 1, M, \dots, N - 1\}$, $\Delta^*(s)$ is continuous and tends to $+\infty$ as s approaches $-aB^{-n}$ and to $-\infty$ as s approaches $-aB^{-(n+1)}$. Therefore, $\Delta^*(s)$ has at least one root in each one of these $N - M + 1$ intervals. Besides that, $\Delta^*(s)$ is also continuous in the interval $(-aB^{-N}, 0]$ and tends to $+\infty$ as s approaches $-aB^{-N}$. When ρ is smaller than 1, as assumed in the proposition, the derivative of $\Delta^*(s)$ is positive at $s = 0$ since

$$\begin{aligned} \left. \frac{d\Delta^*(s)}{ds} \right|_{s=0} &= \left(G^*(s) \frac{dA^*(-s)}{ds} \right) \Big|_{s=0} + \left(A^*(-s) \frac{dG^*(s)}{ds} \right) \Big|_{s=0} \\ &= \frac{1}{\lambda} - \frac{1}{\mu} = \frac{1}{\lambda}(1 - \rho). \end{aligned}$$

Since $\Delta^*(0) = 0$, we conclude that $\Delta^*(s)$ must have at least one root in the interval $(-aB^{-N}, 0)$. We have, thus, found, $N - M + 2$ distinct intervals containing each one at least one root of $\Delta^*(s)$. Reminding that $\Delta^*(s)$ has exactly $N - M + 2$ roots in the left-hand plane, we conclude that a unique root, denoted by σ_n , is contained in each one of the intervals $(-aB^{-n}, -aB^{-(n+1)})$, where $n \in \{M - 1, M, \dots, N - 2, N - 1\}$ and an additional root σ_N is located in the interval $(-aB^{-N}, 0)$.

Appendix C. Proof of proposition 2

In order to prove the first part of the proposition, we set $s = \xi_n$ into $\Delta^*(s)$ and obtain

$$\Delta^*(\xi_n) = -1 + A^*(-\xi_n)G^*(\xi_n). \tag{30}$$

Our goal, now, is to provide asymptotic expansions for $A^*(-\xi_n)$ and $G^*(\xi_n)$ for large and positive values of n . For such values of n , one has $|\xi_n| \ll 1$ and thus $A^*(-\xi_n)$ satisfies the following asymptotic expansion:

$$A^*(-\xi_n) = 1 + \frac{1}{\lambda} \xi_n + o(\xi_n) = 1 - \frac{aB^{-n}}{\lambda} + o(B^{-n}). \tag{31}$$

Regarding $G^*(\xi_n)$, one has

$$\begin{aligned} G^*(\xi_n) &= \frac{(1 - \omega)aB^{-(M-1)}}{-aB^{-n} + cB^{-\alpha n} + aB^{-(M-1)}} \\ &+ \frac{c(1 - \rho)}{\lambda} \sum_{i=M}^N \frac{aB^{-(\alpha+1)i} \exp(-B^{-i})}{-aB^{-n} + cB^{-\alpha n} + aB^{-i}}, \end{aligned} \tag{32}$$

where the constant c is defined as

$$c = \frac{\lambda(B^{\alpha/2} - B^{-\alpha/2})}{(1 - \rho)\Gamma(\alpha + 1)}. \tag{33}$$

We let $N \rightarrow \infty$ and consider large positive values of n , such that $B^{-n} \gg B^{-\alpha n}$. The first term in the right-hand side of equation (32) can then be rewritten in the following way:

$$\begin{aligned} \frac{(1 - \omega)aB^{-(M-1)}}{-aB^{-n} + cB^{-\alpha n} + aB^{-(M-1)}} &= \frac{(1 - \omega)aB^{-(M-1)}}{aB^{-(M-1)}(1 + (-aB^{-n} + cB^{-\alpha n})/(aB^{-(M-1)}))} \\ &= (1 - \omega) \left(1 + \frac{B^{-n}}{B^{-(M-1)}} \right) + o(B^{-n}). \end{aligned} \tag{34}$$

We consider now the second term in the right-hand side of equation (32). We divide the sum appearing in this term into three parts. The first part corresponds to indices i running from M to $n - 1$. We have

$$\begin{aligned} &\sum_{i=M}^{n-1} \frac{aB^{-(\alpha+1)i} \exp(-B^{-i})}{-aB^{-n} + cB^{-\alpha n} + aB^{-i}} \\ &= \sum_{i=M}^{n-1} \frac{aB^{-(\alpha+1)i} \exp(-B^{-i})}{aB^{-i}(1 - (aB^{-n} - cB^{-\alpha n})/(aB^{-i}))} \\ &= \sum_{i=M}^{n-1} B^{-\alpha i} \exp(-B^{-i}) \sum_{k=0}^{\infty} \left(\frac{aB^{-n} - cB^{-\alpha n}}{aB^{-i}} \right)^k \\ &= \sum_{k=0}^{\infty} \left(\frac{aB^{-n} - cB^{-\alpha n}}{a} \right)^k \sum_{i=M}^{n-1} B^{(k-\alpha)i} \exp(-B^{-i}) \\ &= \sum_{i=M}^{n-1} B^{-\alpha i} \exp(-B^{-i}) + \left(\frac{aB^{-n} - cB^{-\alpha n}}{a} \right) \sum_{i=M}^{n-1} B^{(1-\alpha)i} \exp(-B^{-i}) \\ &\quad + \sum_{k=2}^{\infty} \left(\frac{aB^{-n} - cB^{-\alpha n}}{a} \right)^k \sum_{i=M}^{n-1} B^{(k-\alpha)i} \exp(-B^{-i}). \end{aligned} \tag{35}$$

We show now that the third term of equation (35) is in the order of $o(B^{-n})$. We use the notation $\delta_{k\alpha}$ for denoting a function equal to 1 if $k = \alpha$ and to 0 for other values of k . We have

$$\begin{aligned} 0 &< \sum_{k=2}^{\infty} \left(\frac{aB^{-n} - cB^{-\alpha n}}{a} \right)^k \sum_{i=M}^{n-1} B^{(k-\alpha)i} \exp(-B^{-i}) \leq \sum_{k=2}^{\infty} B^{-nk} \sum_{i=M}^{n-1} B^{(k-\alpha)i} \\ &= \sum_{k=2}^{\infty} B^{-nk} \frac{B^{(k-\alpha)M} - B^{(k-\alpha)n}}{1 - B^{(k-\alpha)}} \cdot (1 - \delta_{k\alpha}) + (n - M)B^{-nk} \delta_{k\alpha} \end{aligned}$$

$$\begin{aligned}
 &\leq \max_{k \geq 2} \left(\frac{B^{k-\alpha}}{|1 - B^{k-\alpha}|} \cdot (1 - \delta_{k\alpha}) \right) \sum_{k=2}^{\infty} B^{-nk} |B^{(k-\alpha)(M-1)} - B^{(k-\alpha)(n-1)}| \\
 &\quad + (n - M)B^{-\alpha n} \\
 &\leq \max_{k \geq 2} \left(\frac{B^{k-\alpha}}{|1 - B^{k-\alpha}|} \cdot (1 - \delta_{k\alpha}) \right) \left(\sum_{k=2}^{\infty} B^{-\alpha(M-1)} B^{k(M-1-n)} + B^{-\alpha(n-1)} B^{-k} \right) \\
 &\quad + (n - M)B^{-\alpha n} \\
 &\leq \max_{k \geq 2} \left(\frac{B^{k-\alpha}}{|1 - B^{k-\alpha}|} \cdot (1 - \delta_{k\alpha}) \right) \left(\frac{B^{(2-\alpha)(M-1)}}{1 - B^{M-1} B^{-n}} B^{-2n} + \frac{B^{\alpha-2}}{1 - B} B^{-\alpha n} \right) \\
 &\quad + (n - M)B^{-\alpha n} \\
 &= o(B^{-n}).
 \end{aligned}$$

Note that the expression $(1 - \delta_{k\alpha})B^{k-\alpha}/(1 - B^{k-\alpha})$ is bounded since k is a discrete parameter. We have, thus,

$$\begin{aligned}
 &\sum_{i=M}^{n-1} \frac{aB^{-(\alpha+1)i} \exp(-B^{-i})}{-aB^{-n} + cB^{-\alpha n} + aB^{-i}} \\
 &= \sum_{i=M}^{n-1} B^{-\alpha i} \exp(-B^{-i}) + B^{-n} \sum_{i=M}^{n-1} B^{(1-\alpha)i} \exp(-B^{-i}) + o(B^{-n}). \quad (36)
 \end{aligned}$$

The second part of the sum appearing in the last term of equation (32) corresponds to the index $i = n$ for which

$$\frac{aB^{-(\alpha+1)n} \exp(-B^{-n})}{cB^{-\alpha n}} = \frac{aB^{-n}}{c} \sum_{k=0}^{\infty} \frac{(-B^{-n})^k}{k!} = \frac{aB^{-n}}{c} + o(B^{-n}). \quad (37)$$

The third part of the sum corresponds to indices i larger than or equal to $n + 1$. The contribution of this part is in the order of $o(B^{-n})$ since

$$\begin{aligned}
 &\left| \sum_{i=n+1}^N \frac{aB^{-(\alpha+1)i} \exp(-B^{-i})}{-aB^{-n} + cB^{-\alpha n} + aB^{-i}} \right| \\
 &\leq \frac{\sum_{i=n+1}^N aB^{-(\alpha+1)i}}{\min_{i \geq n+1} |-aB^{-n} + cB^{-\alpha n} + aB^{-i}|} \\
 &= \frac{aB^{-(\alpha+1)(n+1)}}{(1 - B^{-(\alpha+1)}) \cdot |-aB^{-n} + cB^{-\alpha n} + aB^{-(n+1)}|} = o(B^{-n}). \quad (38)
 \end{aligned}$$

Substituting (34), (36)–(38) into (32), and rearranging the terms, we obtain

$$G^*(\xi_n) = \underbrace{(1 - \omega) + \frac{c(1 - \rho)}{\lambda} \sum_{i=M}^{n-1} B^{-\alpha i} \exp(-B^{-i})}_{(I)} + o(B^{-n}) + B^{-n} \underbrace{\left(\frac{(1 - \omega)}{B^{-(M-1)}} + \frac{c(1 - \rho)}{\lambda} \sum_{i=M}^{n-1} B^{(1-\alpha)i} \exp(-B^{-i}) + \frac{a(1 - \rho)}{\lambda} \right)}_{(II)}. \quad (39)$$

We recall the normalization condition which states that

$$(1 - \omega) + \frac{c(1 - \rho)}{\lambda} \sum_{i=M}^N B^{-\alpha i} \exp(-B^{-i}) = 1.$$

Therefore, part (I) in equation (39) can be rewritten as

$$\begin{aligned} & (1 - \omega) + \frac{c(1 - \rho)}{\lambda} \cdot \left(\sum_{i=M}^N B^{-\alpha i} \exp(-B^{-i}) - \sum_{i=n}^N B^{-\alpha i} \exp(-B^{-i}) \right) \\ &= 1 - \frac{c(1 - \rho)}{\lambda} \sum_{i=n}^N B^{-\alpha i} \exp(-B^{-i}) = 1 + o(B^{-n}). \end{aligned} \quad (40)$$

Besides that, the expression

$$\frac{(1 - \omega)}{aB^{-(M-1)}} + \frac{c(1 - \rho)}{a\lambda} \sum_{i=M}^N B^{(1-\alpha)i} \exp(-B^{-i})$$

corresponds to the mean ($1/\mu$) of the pseudo-Pareto distribution. Thus, part (II) in equation (39) can be rewritten as

$$\begin{aligned} & \frac{(1 - \omega)}{B^{-(M-1)}} + \frac{a(1 - \rho)}{\lambda} \\ &+ \frac{c(1 - \rho)}{\lambda} \left(\sum_{i=M}^N B^{(1-\alpha)i} \exp(-B^{-i}) - \sum_{i=n}^N B^{(1-\alpha)i} \exp(-B^{-i}) \right) \\ &= a\mu + \frac{a(1 - \rho)}{\lambda} - \frac{c(1 - \rho)}{\lambda} \sum_{i=n}^N B^{(1-\alpha)i} \exp(-B^{-i}) = \frac{a}{\lambda} + o(1). \end{aligned} \quad (41)$$

Substituting (40) and (41) into (39), we obtain

$$G^*(-\xi_n) = 1 + \frac{aB^{-n}}{\lambda} + o(B^{-n}). \quad (42)$$

We now insert the asymptotic expressions for $A^*(-\xi_n)$ and $G^*(\xi_n)$ given by (31) and (42) into (30) and obtain

$$\Delta^*(\xi_n) = -1 + \left(1 - \frac{aB^{-n}}{\lambda} + o(B^{-n})\right) \left(1 + \frac{aB^{-n}}{\lambda} + o(B^{-n})\right) = o(B^{-n}) \quad (43)$$

which proves the first part of the proposition.

In order to prove the second part of the proposition, we let ε be any constant different from 0 and derive an expression for $\Delta^*[\xi_n + \varepsilon cB^{-\alpha n} + o(B^{-\alpha n})]$. It turns out that all the expressions obtained during the derivation of $\Delta^*(\xi_n)$ remain the same, except for equation (37). Instead, we have

$$\frac{aB^{-(\alpha+1)n} \exp(-B^{-n})}{(1 + \varepsilon)cB^{-\alpha n} + o(B^{-\alpha n})} = \frac{aB^{-n}}{(1 + \varepsilon)c} + o(B^{-n}). \quad (44)$$

We obtain, then,

$$\Delta^*[\xi_n + \varepsilon cB^{-\alpha n} + o(B^{-\alpha n})] = -\frac{\varepsilon ca(1 - \rho)}{\lambda(1 + \varepsilon)} B^{-n} + o(B^{-n}). \quad (45)$$

For sufficiently large values of n , the sign of $\Delta^*[\xi_n + \varepsilon cB^{-\alpha n} + o(B^{-\alpha n})]$ is positive when $\varepsilon < 0$ and negative when $\varepsilon > 0$. From arguments of continuity, it follows that $\sigma_n = \xi_n + o(B^{-\alpha n})$.

Appendix D. Proof of equation (21)

The coefficients ν_n are easily computed by resorting to the residue theorem which gives

$$\nu_n = (\sigma_n + aB^{-n}) \prod_{i=M-1, i \neq n}^N \left(\frac{-\sigma_i}{aB^{-i}}\right) \left(\frac{\sigma_n + aB^{-i}}{\sigma_n - \sigma_i}\right). \quad (46)$$

We derive, now, an explicit asymptotic expression for ν_n as $n \rightarrow \infty$. We substitute the asymptotic expression for σ_n , given by proposition 2, into equation (46) and obtain

$$\nu_n = \{cB^{-\alpha n} + o(B^{-\alpha n})\} \times \prod_{i=M-1, i \neq n}^N \left(\frac{-\sigma_i}{aB^{-i}}\right) \left(\frac{-aB^{-n} + cB^{-\alpha n} + o(B^{-\alpha n}) + aB^{-i}}{-aB^{-n} + cB^{-\alpha n} + o(B^{-\alpha n}) - \sigma_i}\right), \quad (47)$$

where c is defined in the same way as in equation (33). Clearly, equation (21) is proved if one can show that the product appearing in the right-hand side of equation (47) satisfies

$$\prod_{i=M-1, i \neq n}^N \left(\frac{-\sigma_i}{aB^{-i}}\right) \left(\frac{-aB^{-n} + cB^{-\alpha n} + o(B^{-\alpha n}) + aB^{-i}}{-aB^{-n} + cB^{-\alpha n} + o(B^{-\alpha n}) - \sigma_i}\right) = 1 + o(1), \quad (48)$$

as $n \rightarrow \infty$. We let $0 < \varepsilon < \alpha - 1$ and define $\delta = (1 + \varepsilon)/\alpha$ such that $1/\alpha < \delta < 1$. In the sequel, we use the notation $f(n) = O(g(n))$ to mean $\lim_{n \rightarrow \infty} f(n)/g(n) = K$, where $0 < K < \infty$. We consider now the product appearing in the right-hand side of equation (47) and divide it into three parts. The first part corresponds to indices i

running from $M - 1$ to δn (with some abuse of notation, we write δn instead of $\lceil \delta n \rceil$). We have

$$\prod_{i=M-1}^{\delta n} \left(\frac{-\sigma_i}{aB^{-i}} \right) \left(\frac{-aB^{-n} + cB^{-\alpha n} + o(B^{-\alpha n}) + aB^{-i}}{-aB^{-n} + cB^{-\alpha n} + o(B^{-\alpha n}) - \sigma_i} \right)$$

$$= \prod_{i=M-1}^{\delta n} \left(\frac{-\sigma_i}{aB^{-i}} \right) \left(\frac{-aB^{-n} + o(B^{-n}) + aB^{-i}}{-aB^{-n} + o(B^{-n}) - \sigma_i} \right) \tag{49}$$

$$= \prod_{i=M-1}^{\delta n} \left(\frac{-\sigma_i}{aB^{-i}} \right) \left(\frac{-aB^{-n} + o(B^{-n}) + aB^{-i}}{-\sigma_i} \right) \cdot \left(1 - \frac{aB^{-n} + o(B^{-n})}{\sigma_i} \right) \tag{50}$$

$$= \prod_{i=M-1}^{\delta n} \left(\frac{-\sigma_i}{aB^{-i}} \right) \left(\frac{aB^{-i}}{-\sigma_i} \right) \left(1 - \frac{aB^{-n} + o(B^{-n})}{aB^{-i}} \right) \left(1 - \frac{aB^{-n} + o(B^{-n})}{\sigma_i} \right)$$

$$= \prod_{i=M-1}^{\delta n} \left(1 - \frac{aB^{-n} + o(B^{-n})}{aB^{-i}} \right) \left(1 - \frac{aB^{-n} + o(B^{-n})}{\sigma_i} \right)$$

$$= \exp \left[\ln \left(\prod_{i=M-1}^{\delta n} \left(1 - \frac{aB^{-n} + o(B^{-n})}{aB^{-i}} \right) \left(1 - \frac{aB^{-n} + o(B^{-n})}{\sigma_i} \right) \right) \right] \tag{51}$$

$$= \exp \left[- \sum_{i=M-1}^{\delta n} \left(\frac{aB^{-n} + o(B^{-n})}{aB^{-i}} + \frac{aB^{-n} + o(B^{-n})}{\sigma_i} \right) \right] \tag{52}$$

$$= \exp [\mathcal{O}(B^{-(1-\delta)n})] = 1 + o(1). \tag{53}$$

The transitions from (49) to (50) and from (51) to (52) are justified by the fact that for all $i \leq \delta n$ we have $|\sigma_i| > aB^{-(\delta n+1)}$, according to proposition 1, and thus $|\sigma_i| \gg aB^{-n}$.

The second part of the product term corresponds to indices i running from δn to $n - 1$. We note that

$$1 < \prod_{i=\delta n+1}^{n-1} \left(\frac{-aB^{-n} + cB^{-\alpha n} + o(B^{-\alpha n}) + aB^{-i}}{-aB^{-n} + cB^{-\alpha n} + o(B^{-\alpha n}) - \sigma_i} \right)$$

$$= \prod_{i=\delta n+1}^{n-1} \left(1 + \frac{cB^{-\alpha i} + o(B^{-\alpha i})}{-aB^{-n} + o(B^{-n}) - \sigma_i} \right)$$

$$\leq \prod_{i=\delta n+1}^{n-1} \left(1 + \frac{cB^{-\alpha i} + o(B^{-\alpha i})}{-aB^{-n} + aB^{-n+1} + o(B^{-n})} \right)$$

$$= \exp [\mathcal{O}(B^{(-\alpha\delta+1)n})] = \exp [\mathcal{O}(B^{-\varepsilon n})] = 1 + o(1). \tag{54}$$

From (54), we have, therefore, that the second part of the product term is in the order of

$$(1 + o(1)) \prod_{i=\delta n+1}^{n-1} \left(\frac{-\sigma_i}{aB^{-i}} \right). \tag{55}$$

The last part of the product term corresponds to indices i running from $n + 1$ to N . In this case, we have

$$\begin{aligned} 1 &> \prod_{i=n+1}^N \left(\frac{-aB^{-n} + cB^{-\alpha n} + o(B^{-\alpha n}) + aB^{-i}}{-aB^{-n} + cB^{-\alpha n} + o(B^{-\alpha n}) - \sigma_i} \right) \\ &= \prod_{i=n+1}^N \left(1 + \frac{cB^{-\alpha i} + o(B^{-\alpha i})}{-aB^{-n} + o(B^{-n}) - \sigma_i} \right) \\ &\geq \prod_{i=n+1}^N \left(1 + \frac{cB^{-\alpha i} + o(B^{-\alpha i})}{-aB^{-n} + aB^{-n-1} + o(B^{-n})} \right) \\ &= \exp[\mathcal{O}(B^{(-\alpha+1)n})] = 1 + o(1), \end{aligned} \tag{56}$$

and from (56) we obtain that the last part of the product term is in the order of

$$(1 + o(1)) \prod_{i=n+1}^N \left(\frac{-\sigma_i}{aB^{-i}} \right). \tag{57}$$

By multiplying (53) with (55) and (57), we obtain that the product appearing in the right-hand side of equation (47) satisfies

$$\begin{aligned} &\prod_{i=M-1, i \neq n}^N \left(\frac{-\sigma_i}{aB^{-i}} \right) \left(\frac{-aB^{-n} + cB^{-\alpha n} + o(B^{-\alpha n}) + aB^{-i}}{-aB^{-n} + cB^{-\alpha n} + o(B^{-\alpha n}) - \sigma_i} \right) \\ &= (1 + o(1)) \prod_{i=\delta n+1, i \neq n}^N \left(\frac{-\sigma_i}{aB^{-i}} \right). \end{aligned} \tag{58}$$

The product term appearing in the right-hand side of equation (58) is in the order of $1 + o(1)$ as $n \rightarrow \infty$ because

$$\begin{aligned} \prod_{i=\delta n+1, i \neq n}^N \left(\frac{-\sigma_i}{aB^{-i}} \right) &= \prod_{i=\delta n+1, i \neq n}^N \left(\frac{aB^{-i} - cB^{-\alpha i} + o(B^{-\alpha i})}{aB^{-i}} \right) \\ &= \prod_{i=\delta n+1, i \neq n}^N \left(1 - \frac{c}{a} B^{(-\alpha+1)i} + o(B^{(-\alpha+1)i}) \right) \\ &= \exp[\mathcal{O}(B^{(-\alpha+1)\delta n})] = 1 + o(1). \end{aligned} \tag{59}$$

Equation (48) is now simply proven by substituting (59) into (58).

Appendix E. Proof of proposition 3

We let $N \rightarrow \infty$ and take, first, t as a constant. We let $1 < \delta < \alpha$ and divide the sum appearing in the right-hand side of equation (22) into two parts:

$$\Pr(W > t) = \sum_{n=M-1}^{\psi(t)} -\frac{\nu_n}{\sigma_n} e^{\sigma_n t} + \sum_{n=\psi(t)+1}^{\infty} -\frac{\nu_n}{\sigma_n} e^{\sigma_n t}, \tag{60}$$

where $\psi(t) = \ln(t)/(\delta \ln(B))$ (for the simplicity of exposition, we assume that this quantity is integral). The first sum appearing in the right-hand side of (60) can be bounded as follows:

$$0 < \sum_{n=M-1}^{\psi(t)} -\frac{\nu_n}{\sigma_n} e^{\sigma_n t} < \sum_{n=M-1}^{\psi(t)} e^{\sigma_n t} < (\psi(t) - M + 2) \exp[\sigma_{\psi(t)} t], \tag{61}$$

where equation (61) follows from the fact that σ_n is increasing with n . According to proposition 1 we have $\sigma_n > -aB^{-(n+1)}$ for all n and, therefore,

$$\begin{aligned} \psi(t) \exp[\sigma_{\psi(t)} t] &< \psi(t) \exp[-aB^{-\psi(t)-1} t] < \frac{\ln(t) \exp[-aB^{-1} t^{-1/\delta+1}]}{\alpha \ln(B)} \\ &= o(t^{-\alpha+1}). \end{aligned} \tag{62}$$

One concludes that the first sum appearing in the right-hand side of (60) decays to zero faster than $t^{-\alpha+1}$.

The second sum in the right-hand side of (60) can be rewritten as follows:

$$\begin{aligned} &\sum_{n=\psi(t)+1}^{\infty} -\frac{cB^{-\alpha n} + o(B^{-\alpha n})}{-aB^{-n} + cB^{-\alpha n} + o(B^{-\alpha n})} \cdot \exp[(-aB^{-n} + cB^{-\alpha n} + o(B^{-\alpha n}))t] \\ &= \sum_{n=\psi(t)+1}^{\infty} (cB^{-(\alpha+1)n} + o(B^{-(\alpha+1)n})) \cdot \exp[(-aB^{-n} + cB^{-\alpha n} + o(B^{-\alpha n}))t], \end{aligned} \tag{63}$$

where c is defined in the same way as in equation (33). For sufficiently large values of t , equation (63) is bounded from below by

$$\sum_{n=\psi(t)+1}^{\infty} (cB^{-(\alpha+1)n} + o(B^{-(\alpha+1)n})) \cdot \exp[-aB^{-n} t] \tag{64}$$

and from above by

$$\begin{aligned} &\exp[cB^{-\alpha\psi(t)} t] \sum_{n=\psi(t)+1}^{\infty} (cB^{-(\alpha+1)n} + o(B^{-(\alpha+1)n})) \cdot \exp[-aB^{-n} t] \\ &= \exp[t^{-\alpha/\delta+1}] \sum_{n=\psi(t)+1}^{\infty} (cB^{-(\alpha+1)n} + o(B^{-(\alpha+1)n})) \cdot \exp[-aB^{-n} t] \\ &= (1 + O(t^{-\alpha/\delta+1})) \sum_{n=\psi(t)+1}^{\infty} (cB^{-(\alpha+1)n} + o(B^{-(\alpha+1)n})) \cdot \exp[-aB^{-n} t]. \end{aligned} \tag{65}$$

Now, we let $t \rightarrow \infty$ and from equations (64) and (65), we conclude that

$$\sum_{n=\psi(t)+1}^{\infty} -\frac{\nu_n}{\sigma_n} e^{\sigma_n t} \sim \sum_{n=\psi(t)+1}^{\infty} cB^{-(\alpha+1)n} \exp[-aB^{-n}t]. \tag{66}$$

Based on similar arguments as developed in section 2.2, we have for sufficiently large values of t

$$\sum_{n=\psi(t)+1}^{\infty} cB^{-(\alpha+1)n} \exp[-aB^{-n}t] > K_1 t^{-\alpha+1}, \tag{67}$$

where K_1 is some positive constant. Equation (67) is proven by considering its lhs as an approximation of some function $K_2 t^{-\alpha+1}$. This approximation is affected only by the discretization and high-frequency errors. The discretization error is uniformly bounded (see equation (11)). The high-frequency error decays faster than $t^{-\alpha+1}$ since (see equation (26))

$$\begin{aligned} \int_{B^{1/2-\psi(t)}}^{\infty} \frac{s^{\alpha-1} e^{-sat}}{\Gamma(\alpha)} ds &= \int_{B^{1/2-\psi(t)}}^1 \frac{s^{\alpha-1} e^{-sat}}{\Gamma(\alpha)} ds + \int_1^{\infty} \frac{s^{\alpha-1} e^{-sat}}{\Gamma(\alpha)} ds \\ &< \int_{B^{1/2-\psi(t)}}^1 \frac{e^{-sat}}{\Gamma(\alpha)} ds + \int_1^{\infty} \frac{s^{[\alpha]-1} e^{-sat}}{\Gamma(\alpha)} ds \\ &< \frac{\exp[-B^{-\psi(t)}at]}{t\Gamma(\alpha)} + O(e^{-at}) = o(t^{-\alpha+1}). \end{aligned}$$

We obtain, therefore, from equations (60), (62), (66) and (67) that

$$\sum_{n=M-1}^{\infty} -\frac{\nu_n}{\sigma_n} e^{\sigma_n t} \sim \sum_{n=\psi(t)+1}^{\infty} cB^{-(\alpha+1)n} \exp[-aB^{-n}t]. \tag{68}$$

Also, using the same arguments as for the derivation of equation (62) it is easy to show that

$$\sum_{n=0}^{\psi(t)} cB^{-(\alpha+1)n} \exp[-aB^{-n}t] = o(t^{-\alpha+1}). \tag{69}$$

From equations (67)–(69) we finally get

$$\sum_{n=M-1}^N -\frac{\nu_n}{\sigma_n} e^{\sigma_n t} \sim \sum_{n=0}^N cB^{-(\alpha+1)n} \exp[-aB^{-n}t],$$

as $N \rightarrow \infty$ and $t \rightarrow \infty$ which corresponds to the statement of the proposition.

References

[1] J. Abate, G. Choudhury and W. Whitt, Waiting-time tail probabilities in queues with long-tail service-time distributions, *Queueing Systems* 16 (1994) 311–338.

- [2] A. Andersen and B. Nielsen, A Markovian approach for modeling packet traffic with long-range dependence, *IEEE J. Selected Areas Comm.* 16(5) (1998) 749–763.
- [3] P. Billingsley, *Convergence of Probability Measures* (Wiley, New York, 1968).
- [4] O. Boxma, Fluid queues and regular variation, *Performance Evaluation* 27/28 (1996) 699–712.
- [5] O. Boxma and J. Cohen, The M/G/1 queue with heavy-tailed service time distribution, *IEEE J. Selected Areas Comm.* 16(5) (1998) 719–732.
- [6] F. Brichet, J. Roberts, A. Simonian and D. Veitch, Heavy traffic analysis of a storage model with long range dependent on/off sources, *Queueing Systems* 23 (1996) 197–215.
- [7] J. Cohen, Some results on regular variation for distributions in queueing and fluctuation theory, *J. Appl. Probab.* 10 (1973) 343–353.
- [8] J. Cohen, *The Single Server Queue*, 2nd ed. (North-Holland, Amsterdam, 1982).
- [9] A. Feldmann and W. Whitt, Fitting mixture of exponentials to long-tail distributions to analyze network performance models, *Performance Evaluation* 31 (1998) 245–279.
- [10] W. Feller, *An Introduction to Probability Theory and Its Applications*, Vol. II, 2nd ed. (Wiley, New York, 1971).
- [11] M. Greiner, M. Jobmann and C. Klüppelberg, Telecommunication traffic, queueing models, and subexponential distributions, *Queueing Systems* 33(1–3) (1999) 125–152.
- [12] M. Greiner, M. Jobmann and L. Lipsky, The importance of power-tail distributions for telecommunication traffic models, *Oper. Res.* 47(2) (1999) 313–326.
- [13] N. Johnson and S. Kotz, *Continuous Univariate Distributions-1* (Wiley, New York, 1970).
- [14] G. Latouche and V. Ramaswami, *Introduction to Matrix Analytic Methods in Stochastic Modeling* (SIAM, Philadelphia, PA, 1999).
- [15] W. Leland, M. Taqqu, W. Willinger and D. Wilson, On the self-similar nature of Ethernet traffic (extended version), *IEEE/ACM Trans. Networking* 2(1) (1994) 1–15.
- [16] L. Lipsky, *Queueing Theory: A Linear Algebraic Approach* (McMillan, New York, 1992).
- [17] B. Mandelbrot, A fast fractional Gaussian noise generator, *Water Resources Res.* 7(3) (1971) 543–553.
- [18] A. Pakes, On the tails of waiting-time distributions, *J. Appl. Probab.* 12 (1975) 555–564.
- [19] S. Resnick and G. Samorodnitsky, Fluid queues, leaky buckets, on–off processes and teletraffic modeling with highly variable and correlated inputs, in: *Self-Similar Network Traffic and Performance Evaluation*, eds. K. Park and W. Willinger (Wiley, New York, 1998).
- [20] S. Robert and J.-Y. Le Boudec, New models for pseudo self-similar traffic, *Performance Evaluation* 30 (1997) 57–68.
- [21] M. Roughan, D. Veitch and M. Rumsewicz, Computing queue-length distributions for power-law queues, in: *Proc. of Infocom '98*, San Francisco (March 1998).
- [22] D. Starobinski, Quality of service in high speed networks with multiple time-scale traffic, Ph.D. dissertation, Technion – Israel Institute of Technology, Haifa, Israel (May 1999).
- [23] D. Starobinski and M. Sidi, Stochastically bounded burstiness for communication networks, *IEEE Trans. Inform. Theory* 46(1) (2000) 206–212.