



ELSEVIER

Available at
www.ComputerScienceWeb.com
POWERED BY SCIENCE @ DIRECT®

Performance Evaluation 54 (2003) 175–206

**PERFORMANCE
EVALUATION**
An International
Journal

www.elsevier.com/locate/peva

The impact of the service discipline on delay asymptotics

S.C. Borst^{a,b,c}, O.J. Boxma^{a,b,*}, R. Núñez-Queija^{a,b}, A.P. Zwart^b

^a CWI, P.O. Box 94079, 1090 GB Amsterdam, Netherlands

^b Department of Mathematics and Computer Science, Eindhoven University of Technology,
P.O. Box 513, 5600 MB Eindhoven, Netherlands

^c Bell Laboratories, Lucent Technologies, P.O. Box 636, Murray Hill, NJ 07974, USA

Abstract

This paper surveys the $M/G/1$ queue with regularly varying service requirement distribution. It studies the effect of the service discipline on the tail behavior of the waiting-time and/or sojourn-time distribution, demonstrating that different disciplines lead to quite different tail behavior. The orientation of the paper is methodological: We outline four different methods for determining tail behavior, illustrating them for service disciplines like FCFS, Processor Sharing and LCFS.

© 2003 Elsevier Science B.V. All rights reserved.

Keywords: $M/G/1$; Service discipline; Delay asymptotics; Regular variation

1. Introduction

Measurements indicate that traffic in high-speed networks exhibits burstiness on a wide range of time scales, manifesting itself in long-range dependence and self-similarity, see for instance [42,51]. The occurrence of these phenomena is commonly attributed to extreme variability and heavy-tailed characteristics in the underlying activity patterns (connection times, file sizes, scene lengths), see for instance [10,30,56]. This has triggered a lively interest in queueing models with heavy-tailed traffic characteristics.

Although the presence of heavy-tailed traffic characteristics is widely acknowledged, the practical implications for network performance and traffic engineering remain to be fully resolved.

While several studies indicate that small buffer sizes, high levels of aggregation, and flow control algorithms limit the impact on packet-scale buffer dynamics (see, e.g. [5,24]), heavy-tailed traffic characteristics do dramatically affect flow-level delays experienced by users. A particularly interesting aspect is the role of scheduling and priority mechanisms in controlling the latter negative effect.

In a fundamental paper, Anantharam [4] considered a single-server queue fed by a Poisson arrival process of sessions, whose generic length is distributed as some integer random variable T with

* Corresponding author. Present address: Department of Mathematics and Computer Science, Technische Universiteit Eindhoven, P.O. Box 513, 5600 MB Eindhoven, Netherlands. Tel.: +31-40-247-2858; fax: +31-40-246-5995.
E-mail address: boxma@win.tue.nl (O.J. Boxma).

$\mathbf{P}\{T = k\} \sim \alpha k^{-(\alpha+1)} L(k)$, where $1 < \alpha < 2$ and $L(\cdot)$ is a slowly varying function (see Definition 2.1). Each session brings in work at unit rate while it is active. Hence, the work brought in by each arrival is regularly varying and, because $1 < \alpha < 2$, the arrival process of work is long-range dependent, but $\mathbf{E}\{T\} < \infty$. Anantharam shows that, in the steady-state case, for *any* stationary Non-Preemptive service policy, the sojourn-time of a typical session must stochastically dominate a regularly varying random variable having infinite mean. Non-preemption means that once service on a session has begun, it is continued until all the work associated with it has been completed. Anantharam does not make any assumptions as to whether the service policy is work-conserving, or whether the length of a session is known at the time of arrival. In contrast, Anantharam further shows that there also exist causal stationary *preemptive* policies, which do not need information about the session durations at the time of their arrival, for which the sojourn-time of a session is stochastically dominated by a regularly varying random variable with finite mean. The results of Anantharam raise several questions, like (i) are there (preemptive) service disciplines for which the tail of the sojourn-time distribution is not heavier than the tail of the service requirement distribution, and (ii) what is the effect of various well-known scheduling disciplines on the tail behavior of the waiting-time and/or sojourn-time distribution?

A related issue arises when there are *several classes* of customers, which may be treated in different ways by the server (e.g., using fixed priorities, or according to a polling discipline). Then it is important to understand under what conditions, or to what extent, the tail behavior of the service requirements of one class affects the performance of other classes. The above issues have recently been investigated by the present authors and some of their colleagues. This paper summarizes the results. We focus on the classical $M/G/1$ queue and its multi-class generalizations (although some of the recently obtained results allow a general renewal arrival process, or a fluid input).

The orientation of the paper is methodological. After introducing the model and reviewing the main results for various basic disciplines in Section 2, we discuss four different methods for obtaining the tail behavior of waiting-time and/or sojourn-time distributions for $M/G/1$ -type queues with regularly varying service requirement distribution(s): (i) an analytical one, which relies on Tauberian theorems relating the tail behavior of a probability distribution to the behavior of its Laplace–Stieltjes transform near the origin; (ii) a probabilistic one, which exploits a Markov-type inequality, relating an extremely large sojourn (or waiting) time to a single extremely large service requirement; (iii) a probabilistic one, which is based on sample-path arguments which lead to lower and upper bounds for tail probabilities; (iv) a probabilistic one, which is based on explicit (random-sum) representations of the waiting-time distribution, which are applicable to the larger class of subexponential distributions. These four approaches are described in Sections 3–6, respectively. Sections 3, 5 and 6 also discuss the multi-class case. Concluding remarks are given in Section 7. The present paper is an extended version of [16]. In the present version, minor changes have been made in Sections 3 and 4, Section 5 is significantly improved at several points, and Section 6 is new.

2. Model description and main results

In this section, we formally describe the model, introduce some concepts and notation, and give an overview of the main results.

As mentioned earlier, we focus on the $M/G/1$ queue. In this system, customers arrive according to a Poisson process, with rate λ , at a single server who works at unit rate. Their service requirements

B_1, B_2, \dots are independent and identically distributed, with distribution $B(\cdot)$ with mean β and Laplace–Stieltjes transform (LST) $\beta\{\cdot\}$. A generic service requirement is denoted by B . There is no restriction on the number of customers in the system. We assume that the offered traffic load $\rho := \lambda\beta < 1$, so that the system reaches steady-state. We study the steady-state sojourn-time S of a customer, and in some cases also the steady-state waiting-time W until service begins.

Before surveying the tail asymptotics of the waiting-time and/or sojourn-time distributions for various service disciplines, we first introduce some useful notation and terminology. For any two real functions $g(\cdot)$ and $h(\cdot)$, we use the notational convention $g(x) \sim h(x)$ to denote $\lim_{x \rightarrow \infty} g(x)/h(x) = 1$, or equivalently, $g(x) = h(x)(1 + o(1))$ as $x \rightarrow \infty$. For any stochastic variable X with distribution function $F(\cdot)$, with $\mathbf{E}\{X\} < \infty$, denote by $F^r(\cdot)$ the distribution function of the residual lifetime of X , i.e., $F^r(x) = 1/\mathbf{E}\{X\} \int_0^x (1 - F(y)) dy$, and by X^r a stochastic variable with distribution $F^r(\cdot)$.

We focus on the class \mathcal{R} of *regularly varying* distributions (which contains the Pareto distribution). This class is a subset of the class of subexponential distributions [39] as treated in Section 6, which includes for example the lognormal and Weibull distributions as well.

Definition 2.1. A distribution function $F(\cdot)$ on $[0, \infty)$ is called *regularly varying of index $-\nu$* ($F(\cdot) \in \mathcal{R}_{-\nu}$) if

$$1 - F(x) = x^{-\nu} L(x), \quad \nu \geq 0,$$

where $L : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ is a slowly varying function, i.e., $\lim_{x \rightarrow \infty} L(\eta x)/L(x) = 1$, $\eta > 1$.

The class of regularly varying functions was introduced by Karamata [37], and its potential for probability theory was extensively discussed in [33]. A key reference is [12].

In the remainder of this section, we present an overview of the tail asymptotics of the waiting-time and/or sojourn-time distributions in the $M/G/1$ queue for six key disciplines: (i) First-Come-First-Served (FCFS); (ii) Processor Sharing (PS); (iii) Last-Come-First-Served Preemptive-Resume (LCFS-PR); (iv) Last-Come-First-Served Non-Preemptive Priority (LCFS-NP); (v) Foreground-Background Processor Sharing (FBPS); (vi) Shortest-Remaining-Processing-Time-First (SRPTF).

(i) The $M/G/1$ FCFS queue. The next theorem characterizes the tail asymptotics of the distribution of the steady-state waiting-time W for the FCFS service discipline.

Theorem 2.1. *In the case of regular variation, i.e., $\mathbf{P}\{B > \cdot\} \in \mathcal{R}_{-\nu}$,*

$$\mathbf{P}\{W > x\} \sim \frac{\rho}{1 - \rho} \mathbf{P}\{B^r > x\}, \quad x \rightarrow \infty. \quad (1)$$

Remark 2.1. The waiting-time tail in the $M/G/1$ FCFS queue is ‘one degree heavier’ than the service requirement tail, in the regularly varying case. This may be explained by the fact that an arriving customer has a positive probability of arriving during a service time. Its waiting-time is then at least equal to the residual duration of the ongoing service—which is regularly varying of index $1 - \nu$ (cf. [12]).

Theorem 2.1 was first proved by Cohen [25] (who in fact considered the $GI/G/1$ case), and subsequently extended by several authors. In particular, Pakes [49] proved that the relation $\mathbf{P}\{W > x\} \sim \rho/(1 - \rho) \mathbf{P}\{B^r > x\}$ even holds for the larger class of service requirement distributions for which the residual service requirement distribution is subexponential, cf. Section 6.

The fact that the sojourn-time of a customer in the $M/G/1$ FCFS queue equals the sum of its waiting-time and its lighter-tailed service requirement, the two quantities being independent, implies that the tail behavior of the sojourn-time distribution is also given by the right-hand side of (1).

(ii) *The $M/G/1$ PS queue.* The PS (processor sharing) service discipline operates as follows. If there are $n \geq 1$ customers present, then they are all served simultaneously, each at a rate of $1/n$. At the ITC conference in 1997, Roberts raised the question whether the tail of the distribution of the steady-state sojourn-time S_{PS} in the $M/G/1$ PS queue might be just as heavy as the tail of the service requirement distribution. This question was motivated by the following observations: (i) in a PS queue short jobs can overtake long jobs, so the influence of long jobs on the sojourn-time of short jobs is limited, and (ii) the mean sojourn-time in the $M/G/1$ PS queue only involves the *first* moment of the service requirement, whereas in the $M/G/1$ FCFS queue it involves the first moment of the *residual* service requirement, and hence the *second* moment of the service requirement. In fact, if $\mathbf{P}\{B > \cdot\} \in \mathcal{R}_{-\nu}$ with $1 < \nu < 2$, then the second moment of the service requirement does not exist, and neither does the first moment of the waiting-time in the $M/G/1$ FCFS case. Roberts' question can be answered affirmatively, as shown by the next theorem proven in [63].

Theorem 2.2. *If $\mathbf{P}\{B > \cdot\} \in \mathcal{R}_{-\nu}$,*

$$\mathbf{P}\{S_{PS} > x\} \sim \mathbf{P}\{B > (1 - \rho)x\}, \quad x \rightarrow \infty. \quad (2)$$

Apparently, in the $M/G/1$ PS queue the sojourn-time tail is just as heavy as the service requirement tail, which agrees with the observations (i) and (ii).

Remark 2.2. Formula (2) states that the probability that a tagged customer's sojourn-time exceeds the value x is asymptotically (for $x \rightarrow \infty$) equal to the probability that a customer's service requirement exceeds a value $(1 - \rho)x$. This property can be made intuitively plausible as follows: if a customer with an extremely large service requirement is placed in the queue, then the queue remains stable, which heuristically implies that all other customers eventually leave the system. Hence, the average capacity devoted to the service of other customers is approximately equal to ρ . Thus, the average service capacity devoted to the tagged customer is approximately $1 - \rho$. Since the service time of the tagged customer is heavy-tailed, the tagged customer stays in the system long enough to observe this average behavior. Note that Eq. (2) and this intuitive argument are not valid in, e.g., the $M/M/1$ PS queue [15].

(iii) *The $M/G/1$ LCFS-PR queue.* In the LCFS Preemptive-Resume discipline, an arriving customer K is immediately taken into service. However, this service is interrupted when another customer arrives, and it is only resumed when all customers who have arrived after K have left the system.

The fact that no customer has to wait for the completion of a residual service requirement, suggests that the tail of the sojourn-time distribution is just as heavy as the tail of the service requirement distribution. This was indeed proven in [20], using the following observation: the sojourn-time of K has exactly the same distribution as the busy-period of this $M/G/1$ queue. The busy-period obviously has the same distribution for LCFS-PR as for FCFS. The tail behavior of the busy-period distribution in the $M/G/1$ queue has been studied by De Meyer and Teugels [43] for the case of a

regularly varying service requirement distribution. This yields the next theorem (S_{LPR} denoting the steady-state sojourn-time).

Theorem 2.3. *If $\mathbf{P}\{B > \cdot\} \in \mathcal{R}_{-\nu}$,*

$$\mathbf{P}\{S_{\text{LPR}} > x\} \sim \frac{1}{1-\rho} \mathbf{P}\{B > (1-\rho)x\}, \quad x \rightarrow \infty. \quad (3)$$

(iv) *The M/G/1 LCFS-NP queue.* Let W_{LNP} denote the steady-state waiting-time in the M/G/1 LCFS-NP queue. The impossibility of preemption suggests that the tail of W_{LNP} will be determined by the tail of a residual service requirement. Indeed, in this paper we prove the following result, which in fact also holds for the sojourn-time $S_{\text{LNP}} = W_{\text{LNP}} + B$, as is shown in Section 4.

Theorem 2.4. *If $\mathbf{P}\{B > \cdot\} \in \mathcal{R}_{-\nu}$,*

$$\mathbf{P}\{W_{\text{LNP}} > x\} \sim \rho \mathbf{P}\{B^r > (1-\rho)x\}, \quad x \rightarrow \infty. \quad (4)$$

(v) *The M/G/1 FBPS queue.* The Foreground-Background Processor Sharing discipline allocates an equal share of the service capacity to the customers which so far have received the least amount of service, see [38] or [58]. It was proven in [47] (only for the case $1 < \nu < 2$) that the tail of the distribution of the sojourn-time S_{FB} is the same as that for the ordinary PS discipline.

Theorem 2.5. *If $\mathbf{P}\{B > \cdot\} \in \mathcal{R}_{-\nu}$ with $1 < \nu < 2$,*

$$\mathbf{P}\{S_{\text{FB}} > x\} \sim \mathbf{P}\{B > (1-\rho)x\}, \quad x \rightarrow \infty. \quad (5)$$

Although not proven here, it can be shown that the result remains true for $\nu \geq 2$.

(vi) *The M/G/1 SRPTF queue.* With this service discipline the total service capacity is always allocated to the customer(s) with the shortest remaining processing time (Shortest-Remaining-Processing-Time-First). Assuming that $B(x)$ is a continuous function, with probability 1, no two customers in the system have the same remaining service requirement [54]. The service of a customer is preempted when a new customer arrives with a service requirement smaller than the remaining service requirement of the customer being served. The service of the customer that is preempted is resumed as soon as there are no other customers with a smaller amount of work in the system. For the sojourn-time S_{SR} we will prove the following theorem, cf. [47].

Theorem 2.6. *If $\mathbf{P}\{B > \cdot\} \in \mathcal{R}_{-\nu}$ with $1 < \nu < 2$,*

$$\mathbf{P}\{S_{\text{SR}} > x\} \sim \mathbf{P}\{B > (1-\rho)x\}, \quad x \rightarrow \infty. \quad (6)$$

Note that the tail of the service requirement distribution behaves as those of the PS and FBPS disciplines. Again, we remark (without proof) that the result is also valid for $\nu \geq 2$.

In the sequel we prove these theorems using different methods. This serves as an illustration of the various methods and allows us to compare them. Theorems 2.2 and 2.4 are proven in each of Sections 3–6; Theorems 2.1 and 2.3 are proven in Sections 3 and 5, and Theorems 2.5 and 2.6 are proven in Section 4.

3. Transform approach

In this section we outline an LST approach to the study of tails of waiting-time and/or sojourn-time distributions in the $M/G/1$ queue and some of its generalizations. In Section 3.1 we consider the single-class $M/G/1$ queue, with as service discipline either FCFS, PS, LCFS-PR, or LCFS-NP. In Section 3.2 we consider the multi-class $M/G/1$ queue, in which the classes are served according to some scheduling mechanism.

For several of the above-mentioned cases, expressions for the LST of the waiting-time and/or sojourn-time distribution are available in the literature. Such expressions lend themselves to determine the tail behavior of the associated distributions: there exists a very useful relation between the tail behavior of a regularly varying probability distribution and the behavior of its LST near the origin. That relation often enables one to conclude from the form of the LST of the waiting-time and/or sojourn-time distribution, that the distribution itself is regularly varying at infinity. We present this relation in Lemma 3.1.

Let $F(\cdot)$ be the distribution of a non-negative random variable, with LST $\phi\{s\}$ and finite first n moments μ_1, \dots, μ_n (and $\mu_0 = 1$). Define

$$\phi_n\{s\} := (-1)^{n+1} \left[\phi\{s\} - \sum_{j=0}^n \mu_j \frac{(-s)^j}{j!} \right].$$

Lemma 3.1. *Let $n < \nu < n + 1, C \geq 0$. The following statements are equivalent:*

$$\begin{aligned} \phi_n\{s\} &= (C + o(1))s^\nu L\left(\frac{1}{s}\right), \quad s \downarrow 0, \quad s \text{ real}, \\ 1 - F(x) &= (C + o(1)) \frac{(-1)^n}{\Gamma(1 - \nu)} x^{-\nu} L(x), \quad x \rightarrow \infty. \end{aligned}$$

The case $C > 0$ is due to Bingham and Doney [11]. The case $C = 0$ was first obtained by Vincent Dumas, and is treated in [23, Lemma 2.2]. The case with ν integer-valued is more complicated; see Theorem 8.1.6 and Chapter 3 of [12].

3.1. The single-class case

(i) *The $M/G/1$ FCFS queue.* In the $M/G/1$ FCFS queue, the LST of the steady-state waiting-time distribution is given by the Pollaczek–Khintchine formula [26]:

$$\mathbf{E}\{e^{-sW}\} = \frac{1 - \rho}{1 - \rho\beta^r\{s\}}, \quad \text{Re } s \geq 0, \tag{7}$$

where $\beta^r\{s\} = (1 - \beta\{s\})/\beta s$ is the LST of the residual service requirement distribution $B^r(\cdot)$. A Karamata theorem (cf. Section 1.5 of [12]) implies that if $\mathbf{P}\{B > \cdot\} \in \mathcal{R}_{-\nu}$, then the integrated tail $\mathbf{P}\{B^r > \cdot\} \in \mathcal{R}_{1-\nu}$. More precisely, if

$$\mathbf{P}\{B > x\} \sim x^{-\nu} L(x), \quad \nu > 1, \quad x \rightarrow \infty, \tag{8}$$

then

$$\mathbf{P}\{B^r > x\} = \frac{1}{\beta} \int_x^\infty \mathbf{P}\{B > y\} dy \sim \frac{1}{(v-1)\beta} x^{1-v} L(x), \quad x \rightarrow \infty.$$

We now demonstrate how the following statement, which implies [Theorem 2.1](#), is easily obtained from the LST expression (7) and [Lemma 3.1](#). For $1 < v < 2$, $x \rightarrow \infty$,

$$\mathbf{P}\{B > x\} \sim x^{-v} L(x) \Leftrightarrow \mathbf{P}\{W > x\} \sim \frac{\rho}{1-\rho} \frac{1}{(v-1)\beta} x^{1-v} L(x). \tag{9}$$

It follows from (8) and [Lemma 3.1](#) that

$$1 - \beta^r\{s\} = 1 - \frac{1 - \beta\{s\}}{\beta s} = - \left(\frac{\Gamma(1-v)}{\beta} + o(1) \right) s^{\nu-1} L\left(\frac{1}{s}\right), \quad s \downarrow 0. \tag{10}$$

Combining this result with (7) yields:

$$1 - \mathbf{E}\{e^{-sW}\} = \frac{\rho(1 - \beta^r\{s\})}{1 - \rho\beta^r\{s\}} \sim - \frac{\rho}{1-\rho} \frac{\Gamma(1-v)}{\beta} s^{\nu-1} L\left(\frac{1}{s}\right), \quad s \downarrow 0.$$

Another application of [Lemma 3.1](#) gives the \Rightarrow part of (9). The reverse part is obtained in a similar way. A similar approach can be followed for non-integer values of $v > 2$; we ignore the subtleties required in applying [Lemma 3.1](#) for integer values of v .

(ii) *The M/G/1 PS queue.* [Theorem 2.2](#) indicates that, contrary to the FCFS case, the sojourn-time tail in the M/G/1 PS queue is just as heavy as the service requirement tail. We now sketch the proof in [63], which is based on the application of [Lemma 3.1](#) to an explicit expression of the sojourn-time LST.

There are several expressions known for the LST of the sojourn-time, cf. [48,55,57], but they contain contour integrals which are inversion formulas of Laplace transforms. Starting-point in [63] is an expression in [48] for the conditional LST of a customer’s sojourn-time $S_{PS}(\tau)$, given that his service requirement is τ : for $\text{Re } s \geq 0, \tau \geq 0$,

$$\mathbf{E}\{e^{-sS_{PS}(\tau)}\} = \frac{1 - \rho}{(1 - \rho)H_1(s, \tau) + sH_2(s, \tau)},$$

where the functions $H_1(s, \tau)$ and $H_2(s, \tau)$ are given by their LST w.r.t. τ

$$\int_0^\infty e^{-x\tau} dH_1(s, \tau) = \frac{x - \lambda(1 - \beta\{x\})}{x - s - \lambda(1 - \beta\{x\})}, \quad \text{Re } x > 0,$$

$$\int_0^\infty e^{-x\tau} dH_2(s, \tau) = \frac{\rho x - \lambda(1 - \beta\{x\})}{x(x - s - \lambda(1 - \beta\{x\}))}, \quad \text{Re } x > 0.$$

It follows from these relations that, for $\text{Re } s \geq 0$ and $\text{Re } x > 0$,

$$\int_0^\infty e^{-x\tau} d[\mathbf{E}\{e^{-sS_{PS}(\tau)}\}]^{-1} = 1 + \frac{1}{1-\rho} \frac{s}{x} \frac{1}{1 - s\mathbf{E}\{e^{-xW}\}/(x(1-\rho))}, \tag{11}$$

where W denotes the steady-state waiting-time in the M/G/1 FCFS queue (we will denote its distribution by $W(\cdot)$). Formula (11) implies (see [63]) that

$$\mathbf{E}\{e^{-sS_{PS}(\tau)}\} = \left[\sum_{k=0}^\infty \frac{s^k}{k!} \alpha_k(\tau) \right]^{-1}, \tag{12}$$

with $\alpha_0(\tau) := 1$, $\alpha_1(\tau) := \tau/(1 - \rho)$, and for $k \geq 2$,

$$\alpha_k(\tau) := \frac{k}{(1 - \rho)^k} \int_{x=0}^{\tau} (\tau - x)^{k-1} W^{(k-1)*}(x) dx.$$

In Corollary 3.2 of [63], Eq. (12) is shown to imply that the k th moment of the sojourn-time in the $M/G/1$ PS queue is finite iff the k th moment of the service requirement is finite. But Eq. (12) is also suitable for applying Lemma 3.1. With S_{PS} the steady-state sojourn-time in the $M/G/1$ PS queue, and using the fact that $\mathbf{E}\{e^{-sS_{PS}}\} = \int_0^\infty \mathbf{E}\{e^{-sS(\tau)}\} dB(\tau)$, it can be shown [63] that, for $1 < \nu < 2$,

$$\mathbf{E}\{e^{-sS_{PS}}\} - \beta \left\{ \frac{s}{1 - \rho} \right\} = o \left(s^\nu L \left(\frac{1}{s} \right) \right), \quad s \downarrow 0, \quad s \text{ real.}$$

One can now apply Lemma 3.1. Using the well-known fact that $\mathbf{E}\{S_{PS}\} = \beta/(1 - \rho)$, it is seen that Theorem 2.2 holds for $1 < \nu < 2$ (and via a similar approach it is shown in [63] that this holds for all non-integer $\nu > 1$). In fact, a two-way application of Lemma 3.1 yields (cf. [63]): for non-integer $\nu > 1$, $x \rightarrow \infty$,

$$\mathbf{P}\{B > x\} \sim x^{-\nu} L(x) \Leftrightarrow \mathbf{P}\{S_{PS} > x\} \sim \frac{1}{(1 - \rho)^\nu} x^{-\nu} L(x).$$

(iii) *The $M/G/1$ LCFS-PR queue.* As observed in Section 2, the sojourn-time in the $M/G/1$ LCFS-PR queue has the same distribution as the busy-period in the $M/G/1$ queue. De Meyer and Teugels [43] have studied the tail of the latter distribution in the case of a regularly varying service requirement distribution. Their starting-point is the fact that the LST $\mu\{s\}$ of the steady-state busy-period length P is the unique solution of the equation

$$\mu\{s\} = \beta\{s + \lambda(1 - \mu\{s\})\} \tag{13}$$

with $|\mu\{s\}| \leq 1$ for $\text{Re } s \geq 0$. They apply Lemma 3.1 to show the following equivalence: for $\nu > 1$, $x \rightarrow \infty$,

$$\mathbf{P}\{B > x\} \sim x^{-\nu} L(x) \Leftrightarrow \mathbf{P}\{P > x\} \sim \frac{1}{(1 - \rho)^{\nu+1}} x^{-\nu} L(x). \tag{14}$$

Hence, the tail of the busy-period distribution is just as heavy as that of the service requirement distribution. Theorem 2.3 immediately follows from (14).

(iv) *The $M/G/1$ LCFS-NP queue.* Let W_{LNP} denote the steady-state waiting-time in the $M/G/1$ LCFS-NP queue. The following is observed in [26, p. 431]. If an arriving customer in the $M/G/1$ LCFS-NP queue meets a customer in service with a residual service requirement w , then its waiting-time distribution is that of a busy-period with a special first service requirement w . That residual service requirement has distribution $B^r(\cdot)$ with LST $\beta^r\{s\}$ as introduced in the beginning of this section [26, p. 432]. It is now readily seen (cf., e.g., [27, p. 299]) that

$$\mathbf{E}\{e^{-sW_{LNP}}\} = 1 - \rho + \rho\beta^r\{\delta\{s\}\}, \quad \text{Re } s \geq 0, \tag{15}$$

with $\delta\{s\}$ the unique zero in $\text{Re } s \geq 0$ of $\lambda(1 - \beta\{w\}) - w + s$, $\text{Re } w \geq 0$. In fact, cf. (13), $\delta\{s\} = s + \lambda(1 - \mu\{s\})$. In combination with (15), this gives an alternative derivation of Eq. (III.3.10) of [26]:

$$\mathbf{E}\{e^{-sW_{LNP}}\} = 1 - \rho + \frac{\rho}{\beta} \frac{1 - \mu\{s\}}{s + \lambda(1 - \mu\{s\})}, \quad \text{Re } s \geq 0. \tag{16}$$

Using Lemma 3.1, we can now easily verify that the tail of W_{LNP} is regularly varying of degree one heavier than the tail of the service requirement (as may be expected in view of the possibility of having to wait at least a residual service requirement). If (8) and hence also (10) hold, then it follows from (14) and Lemma 3.1 that

$$1 - \mu\{s\} - \frac{\beta}{1 - \rho}s \sim -\frac{\Gamma(1 - \nu)}{(1 - \rho)^{\nu+1}}s^\nu L\left(\frac{1}{s}\right), \quad s \downarrow 0, \tag{17}$$

and therefore

$$1 - \mathbf{E}\{e^{-sW_{LNP}}\} \sim -\frac{\lambda\Gamma(1 - \nu)}{(1 - \rho)^{\nu-1}}s^{\nu-1}L\left(\frac{1}{s}\right), \quad s \downarrow 0. \tag{18}$$

On the other hand, starting from (18) and using (16), one gets (17). Application of Lemma 3.1 and (14) now yields: for $\nu > 1, x \rightarrow \infty$,

$$\mathbf{P}\{B > x\} \sim x^{-\nu}L(x) \Leftrightarrow \mathbf{P}\{W_{LNP} > x\} \sim \frac{\lambda}{(\nu - 1)(1 - \rho)^{\nu-1}}x^{1-\nu}L(x).$$

Both relations imply Theorem 2.4.

3.2. The multi-class case

In this section we consider the $M/G/1$ queue with K classes of customers. We study several of the most important service disciplines, rules that specify at any time which class of customers is being served. We are interested in the question under what conditions, or to what extent, the tail behavior of the service requirements of one class affects the performance of other classes.

The notation is as introduced in Section 2, but quantities relating to class- i customers receive an index i . Hence, class- i customers arrive according to a Poisson process with rate λ_i , and their service requirements have distribution $B_i(\cdot)$ with mean β_i ; $\rho_i := \lambda_i\beta_i$ and $\rho := \sum_{i=1}^K \rho_i$.

(i) *Fixed priorities: Non-Preemptive priority.* Assume that there are only two priority classes, class 1 having Non-Preemptive priority over class 2. Cohen [26, Section III.3.8] gives the following expressions for the LST of the distribution of the steady-state waiting-time W_1 of class-1 customers:

$$\mathbf{E}\{e^{-sW_1}\} = \frac{1 - \rho + \rho_2\beta_2^r\{s\}}{1 - \rho_1\beta_1^r\{s\}}, \quad \text{Re } s \geq 0, \quad \rho < 1, \tag{19}$$

$$\mathbf{E}\{e^{-sW_1}\} = \frac{(1 - \rho_1)\beta_2^r\{s\}}{1 - \rho_1\beta_1^r\{s\}}, \quad \text{Re } s \geq 0, \quad \rho_1 < 1, \quad \rho \geq 1. \tag{20}$$

In both cases, Lemma 3.1 can readily be applied to determine the tail behavior of the waiting-time distribution. Actually, this is one of the rare cases in which the LST can be easily inverted. For $\rho < 1$ this gives ($B_{1,i}^r$ has the residual service requirement distribution $B_1^r(\cdot)$, and B_2^r has the residual service requirement distribution $B_2^r(\cdot)$), with $\stackrel{d}{=}$ denoting equality in distribution,

$$W_1 \stackrel{d}{=} B_{1,1}^r + \dots + B_{1,N}^r + Z,$$

where N is geometrically distributed with parameter ρ_1 while Z is zero with probability $(1 - \rho)/(1 - \rho_1)$ and $Z = B_2^r$ with probability $\rho_2/(1 - \rho_1)$. For $\rho_1 < 1$ but $\rho \geq 1$, inversion of the LST in (20) yields:

$$W_1 \stackrel{d}{=} B_{1,1}^r + \dots + B_{1,N}^r + B_2^r.$$

These results imply the following. If the service requirement distribution with the heaviest tail is regularly varying at infinity of index $-\nu$, then the waiting-time distribution of the high-priority customers is regularly varying at infinity of index $1 - \nu$. More specifically: if the heaviest tail belongs to class 1, then the waiting-time tail of class-1 customers is as if no class 2 exists. If the heaviest tail belongs to class 2, then the waiting-time tail of class-1 customers behaves like the tail of a residual service requirement of class 2 if $\rho_1 < 1$ and $\rho \geq 1$, and like that tail multiplied by the factor $\rho_2/(1 - \rho_1)$ if $\rho < 1$.

For class 2 the following result has been proven in [21]. If the service requirement distribution with the heaviest tail is regularly varying at infinity of index $-\nu$, then the waiting-time distribution of the low-priority customers is regularly varying at infinity of index $1 - \nu$. This is proven by exploiting a representation for the LST of that waiting-time distribution, as given by Abate and Whitt [2], and then using Lemma 3.1. The result is not surprising, when one realizes that a low-priority customer may have to wait for a residual service requirement of either class. See [21] for more details. Alternative approaches to this model may be found in Sections 5 and 6.

(ii) *Fixed priorities: Preemptive-Resume priority.* First assume that there are only two priority classes, class 1 having Preemptive-Resume priority over class 2. Clearly, class-1 customers are not affected by class-2 customers, so the results of Section 3.1 (for FCFS) apply to class 1. The waiting-time distribution of the low-priority customers *until the start of the—possibly interrupted—service* is the same as in the Non-Preemptive case. Those possible interruptions consist of full service requirements of high-priority customers, and in the regularly varying case these are less heavy than *residual* service requirements of those customers. Hence, in the scenario of regular variation, the tail behavior of low-priority customers is the same as in the Non-Preemptive case.

If there are $K > 2$ classes, then in studying class j one may aggregate classes $1, \dots, j - 1$ into one high-priority class w.r.t. class j , while the existence of classes $j + 1, \dots, K$ is irrelevant for class j .

(iii) *Polling.* Deng [31] has considered the extension of the two-class Non-Preemptive priority model to the case in which the server requires a switchover time to move from one class of customers to the other. She proves: if the service requirement distribution *or the switchover-time distribution* with the heaviest tail is regularly varying at infinity of index $-\nu$, then the waiting-time distributions of both classes are regularly varying at infinity of index $1 - \nu$. Again the key of the derivation is an explicit expression for the LST of the waiting-time distributions, in combination with Lemma 3.1.

The above two-class model may also be viewed as a *polling* model with two queues Q_1, Q_2 and a server who alternately visits both queues, serving Q_1 exhaustively (i.e., until it is empty) and applying the 1-limited service discipline at Q_2 (i.e., serving one customer, if there is one, and then moving on to the other queue). In [22] a polling model with K queues has been studied, with the exhaustive or gated service discipline being employed at the various queues. In a similar way, the same conclusions as above have been obtained.

(iv) *Processor sharing with several customer classes.* In the multi-class disciplines that were discussed above, the worst tail behavior of any class determined the waiting-time tail behavior of all classes (except for high-priority customers in the case of Preemptive-Resume priority). Processor sharing turns out to be better capable of protecting customer classes from the bad behavior of other classes. Zwart [59] showed that the sojourn-time distribution of a class- i customer is regularly varying of index $-\nu_i$ iff the service requirement distribution of that class is regularly varying of index $-\nu_i$, *regardless* of the service requirement distributions of the other classes. His method again relied on Lemma 3.1.

(v) *Generalized processor sharing.* The generalized processor sharing (GPS) discipline operates as follows [50]. Customer class i is assigned a weight ϕ_i , $i = 1, \dots, K$, with $\sum_{i=1}^K \phi_i = 1$. If customers of

all classes are present, then one customer from each class is served simultaneously (processor sharing), a class- i customer receiving a fraction ϕ_i of the server capacity. If only some of the classes are present, then the service capacity is shared in proportion to the weights ϕ_i among the head-of-the-line customers of those classes.

GPS-based scheduling algorithms, such as Weighted Fair Queueing, play a major role in achieving differentiated quality-of-service in integrated-services networks. Hence, it is important to study the extent to which GPS is capable of protecting one class of customers from the adverse effects of bad traffic characteristics of other classes. Unfortunately, the queueing analysis of GPS is very difficult. A slightly more general model for $K = 2$ is the model with two parallel $M/G/1$ queues with service speeds depending on whether the other queue is empty or not. For general service requirement distributions, the joint distribution of the amounts of work of both classes has been obtained in [27] by solving a Wiener–Hopf problem (see [32,40] for the case of exponential service requirement distributions). The results of [27] have been exploited in [13,17,18]. In those papers, service requirements at Q_1 are either exponential or regularly varying; at Q_2 they are regularly varying. Whether the service requirement tail behavior at Q_2 affects the workload tail at Q_1 is shown to depend crucially on whether or not Q_1 is able to handle all its offered work by working at the low speed that occurs while Q_2 is non-empty (i.e., whether or not $\rho_1 < \phi_1$). The method employed in [13,17,18] starts from a, complicated, expression for the workload LST. In some cases Lemma 3.1 is applicable, but in other cases an extension of this lemma must be used. For $K \geq 3$ coupled queues, respectively, for GPS with $K \geq 3$ classes, no explicit results are known. However, the sample-path techniques discussed in Section 5 have proven useful in obtaining tail asymptotics for an arbitrary number of classes [14].

4. Tail equivalence via conditional moments

With heavy-tailed distributions, it is often the case that large occurrences of the variable of interest (e.g., a customer's waiting-time or sojourn-time) are essentially caused by a *single* large occurrence of one input variable (e.g., a service requirement). In this section we describe a generic approach that may be used to prove that the tails of the distributions of the *causal variable* and the *resultant variable* are *equally heavy*. Specifically, we say that two non-negative random variables X and Y have equally heavy-tailed distributions if $\mathbf{P}\{Y > \bar{g}x\} \sim \mathbf{P}\{X > x\}$ for some constant $\bar{g} > 0$. In the examples below, a customer's own service requirement (denoted by B) or the *residual* service requirement of some other customer (denoted by B') will play the role of the causal variable X and the customer's sojourn-time (denoted by S) that of the resultant Y . In order to explicitly express the dependence of the sojourn-time on the (residual) service requirement, we use $S(\tau)$ to denote a customer's sojourn-time *given* that the (residual) service requirement equals τ . Consequently, we may alternatively write $S(B)$ or $S(B')$ (depending on the causal variable) for the unconditional sojourn-time S .

Theorem 4.1 relates the tails of the distributions of S and the causal variable X (later replaced with either B or B'). We make two assumptions: one regarding the distribution of the causal variable and one regarding $S(\tau)$.

Assumption 4.1. $\mathbf{P}\{X > \cdot\} \in \mathcal{R}_{-\alpha}$ for some $\alpha > 0$.

This assumption can be relaxed to distributions of *intermediate regular variation*, a class introduced by Cline [29], without invalidating Theorem 4.1, see [47].

Assumption 4.2. The following three conditions are satisfied:

- (a) $\mathbf{E}\{S(\tau)\} \sim \bar{g}\tau$, for some $\bar{g} > 0$;
- (b) With α as in Assumption 4.1, there exists $\kappa > \alpha$ such that

$$\mathbf{P}\{S(\tau) - \mathbf{E}\{S(\tau)\} > t\} \leq \frac{h(\tau)}{t^\kappa}$$

with $h(\tau) = o(\tau^{\kappa-\delta})$, $\tau \rightarrow \infty$, for some $\delta > 0$;

- (c) $S(\tau)$ is stochastically increasing in $\tau \geq 0$, i.e., for all $t \geq 0$, the probability $\mathbf{P}\{S(\tau) > t\}$ is non-decreasing in $\tau \geq 0$.

Theorem 4.1. Suppose Assumptions 4.1 and 4.2 are satisfied. Then the tails of the distributions of the random variables X and $S(X)$ are equally heavy in the sense that:

$$\mathbf{P}\{S(X) > \bar{g}x\} \sim \mathbf{P}\{X > x\}.$$

In particular, the distribution of $S(X)$ is also regularly varying with the same index $-\alpha$ as that of X .

Proof. We only give a sketch of the proof and refer to [47] for details. The proof consists of two parts. For the first part we write, with $\varepsilon > 0$,

$$\mathbf{P}\{S(X) > \bar{g}x\} \leq \mathbf{P}\{S(X) > \bar{g}x; X \leq x(1 - \varepsilon)\} + \mathbf{P}\{X > x(1 - \varepsilon)\}. \tag{21}$$

By conditioning on X and integrating over the distribution of X , it can be shown (using Assumptions 4.1 and 4.2) that

$$\mathbf{P}\{S(X) > \bar{g}x; X \leq x(1 - \varepsilon)\} = o(\mathbf{P}\{X > x(1 - \varepsilon)\}), \quad x \rightarrow \infty. \tag{22}$$

Hence, we may neglect the first term on the right-hand side of (21) and write

$$\limsup_{x \rightarrow \infty} \frac{\mathbf{P}\{S(X) > \bar{g}x\}}{\mathbf{P}\{X > x\}} \leq \limsup_{x \rightarrow \infty} \frac{\mathbf{P}\{X > x(1 - \varepsilon)\}}{\mathbf{P}\{X > x\}} = (1 - \varepsilon)^{-\alpha}.$$

Letting $\varepsilon \downarrow 0$, the right-hand side tends to 1.

For the second part of the proof we write, for $\varepsilon > 0$,

$$\mathbf{P}\{S(X) > \bar{g}x\} \geq \mathbf{P}\{S(X) > \bar{g}x; X > x(1 + \varepsilon)\}.$$

By conditioning again on X , it can be shown that

$$\lim_{x \rightarrow \infty} \frac{\mathbf{P}\{S(X) > \bar{g}x; X > x(1 + \varepsilon)\}}{\mathbf{P}\{X > x(1 + \varepsilon)\}} = 1. \tag{23}$$

Hence,

$$\liminf_{x \rightarrow \infty} \frac{\mathbf{P}\{S(X) > \bar{g}x\}}{\mathbf{P}\{X > x\}} \geq \liminf_{x \rightarrow \infty} \frac{\mathbf{P}\{X > x(1 + \varepsilon)\}}{\mathbf{P}\{X > x\}} = (1 + \varepsilon)^{-\alpha}.$$

Again, the right-hand side tends to 1 as $\varepsilon \downarrow 0$. □

Remark 4.1. Formulas (22) and (23) in fact allow us to prove the stronger statement that essentially X and $S(X)$ can only ‘simultaneously exceed’ the values x and $\bar{g}x$, respectively, for $x \rightarrow \infty$.

We will employ [Theorem 4.1](#) to show for several queueing models that the tail of the sojourn-time distribution is as heavy as that of the (residual) service requirement distribution. Assuming that the service requirement distribution is regularly varying, it suffices to verify that $S(\tau)$, the sojourn-time conditioned on the (residual) service requirement, satisfies [Assumption 4.2](#). Parts (a) and (c) of [Assumption 4.2](#) are often not hard to verify. We will use the following variant of Markov’s inequality to verify part (b):

$$\mathbf{P}\{S(\tau) > t\} \leq \frac{\mathbf{E}\{S(\tau)^\kappa\} - (\mathbf{E}\{S(\tau)\})^\kappa}{(t - \mathbf{E}\{S(\tau)\})^\kappa}, \quad \tau \geq 0, \quad t > \mathbf{E}\{S(\tau)\}, \tag{24}$$

where $\kappa \geq 2$. In [\[47\]](#) Markov’s inequality itself was used, but for the analysis of the $M/G/1$ LCFS-NP below, the form of [\(24\)](#) is more convenient. To see that this inequality holds, let $c(y)$, $y \geq 0$, be a convex function with a convex derivative $c'(y)$ and $c(0) = c'(0) = 0$. If Y is a non-negative random variable and $t \geq \mathbf{E}\{Y\}$, then

$$c(Y) - c(\mathbf{E}\{Y\}) - c'(\mathbf{E}\{Y\})(Y - \mathbf{E}\{Y\}) \geq \mathbf{1}_{\{Y>t\}}c(t - \mathbf{E}\{Y\}),$$

where $\mathbf{1}_{\{\cdot\}}$ denotes the indicator function. Taking expectations with respect to the distribution of Y , we obtain

$$\mathbf{E}\{c(Y)\} - c(\mathbf{E}\{Y\}) \geq \mathbf{P}\{Y > t\}c(t - \mathbf{E}\{Y\}).$$

Choosing $c(y) = y^\kappa$, with $\kappa \geq 2$, leads to the desired result.

The strength of the method described here is that it does not rely on the availability of the LST for the sojourn-time distribution. In particular, the method’s flexibility was demonstrated in [\[45\]](#) where it was employed in the analysis of an $M/G/1$ PS queue with random service interruptions (for that model even basic performance measures such as mean queue length are not available). A limitation of the method is that it relies on the fact that an extreme occurrence of the performance measure of interest (e.g., the sojourn-time) is essentially caused by the occurrence of a *single* extreme input variable.

(i) *The M/G/1 PS queue.* Consider again the $M/G/1$ PS queue described in [Section 2](#). In this section $S_{PS}(\tau)$ will stand for the sojourn-time of a customer with service requirement $\tau \geq 0$, arriving when the system has reached steady-state. As before, the unconditional sojourn-time will be denoted by S_{PS} , i.e., $S_{PS} = S_{PS}(B)$, where the random variable B stands for the customer’s service requirement. We first list some known results for the moments of $S_{PS}(\tau)$. Then we will use these to verify [Assumption 4.2](#) and subsequently apply [Theorem 4.1](#).

It is well-known [\[38,53\]](#) that the mean of the conditional sojourn-time is proportional to the service requirement

$$\mathbf{E}\{S_{PS}(\tau)\} = \frac{\tau}{1 - \rho}. \tag{25}$$

The variance of $S_{PS}(\tau)$ is given by

$$\text{Var}\{S_{PS}(\tau)\} = \frac{2}{(1 - \rho)^2} \int_{u=0}^{\tau} (\tau - u)\mathbf{P}\{W > u\} du, \tag{26}$$

cf. [\[57\]](#). As before, W is distributed as the steady-state waiting-time in the $M/G/1$ queue. When the second moment of the service requirement distribution is finite, we have for $k = 2, 3, \dots$, cf. [\[63\]](#),

$$\mathbf{E}\{S_{PS}(\tau)^k\} = \left(\frac{\tau}{1 - \rho}\right)^k + \frac{\lambda k(k - 1)\mathbf{E}\{B^2\}}{2(1 - \rho)^{k+1}}\tau^{k-1} + o(\tau^{k-1}), \quad \tau \rightarrow \infty. \tag{27}$$

In the literature these results have mostly been obtained from expressions for the LST of $S_{PS}(\tau)$. However, (25)–(27) can be obtained directly from a set of differential equations instead of deriving the LST of $S_{PS}(\tau)$, see [57, Remark 3] for an outline of how this can be done for the variance of $S_{PS}(\tau)$, and [45,46] for higher moments. Later in this section we use similar ideas to derive differential equations for the moments of the conditional sojourn-time in the $M/G/1$ LCFS-NP.

We now provide a new proof of Theorem 2.2. As before, we assume that $B(\cdot) \in \mathcal{R}_{-\nu}$. We further require that $\nu \neq 2$; in [47] it is indicated how this technical condition can be avoided. Focusing on the sojourn-time of a particular customer, its own service requirement B will act as the causal variable X , i.e., in the light of Assumption 4.1 we choose $\alpha = \nu$. We now verify that Assumption 4.2 is automatically satisfied. First we note that the monotonicity of $\mathbf{P}\{S_{PS}(\tau) > t\}$ in τ , the last condition in Assumption 4.2, is easily verified using a sample-path argument: comparing the sojourn-times of two customers, for the same sequences of inter-arrival times and service requirements of other customers, it follows immediately that the one requiring the smaller amount of service leaves before the one with the larger service requirement. As a consequence of (25), we also have that Condition (a) of Assumption 4.2 holds with $\bar{g} = 1/(1 - \rho)$.

We now focus on Condition (b) and first consider the case that $\nu > 2$, ensuring that $\mathbf{E}\{B^2\} < \infty$. Choose any integer $\kappa > \nu$ and use (27) to conclude that Condition (b) is satisfied for any $\delta \in (0, 1)$. Hence, Theorem 4.1 can be applied.

In the case that $1 < \nu < 2$, it follows from (26) that $\text{Var}\{S_{PS}(\tau)\} = o(\tau^{3-\nu+\varepsilon})$ for all $\varepsilon > 0$. Thus, Assumption 4.2 is satisfied (with $\kappa = 2$ and $0 < \delta < \nu - 1$) and, hence, Theorem 4.1 can again be applied. In both cases we conclude that, cf. Theorem 2.2,

$$\mathbf{P}\left\{S_{PS} > \frac{x}{1 - \rho}\right\} \sim \mathbf{P}\{B > x\}.$$

(ii) *The M/G/1 LCFS-NP queue.* In the LCFS-NP case we focus on the sojourn-time, $S_{LNP}(\tau)$, of a tagged customer entering the system when the remaining service requirement of the customer in service equals τ . Because of the service discipline, if there are any customers in the queue, these are overtaken by the new customer and they will have no influence on $S_{LNP}(\tau)$, which we may write as

$$S_{LNP}(\tau) = \tau + \sum_{n=1}^{N(\tau)} P_n + B,$$

where, by convention, we set the empty sum equal to 0. $N(\tau)$ denotes the number of customers that enter the system during the remaining service requirement τ of the customer in service. $P_n, n = 1, 2, \dots$, is an i.i.d. sequence having the distribution of the busy-period in the $M/G/1$ queue. Indeed, all customers that enter during the time τ overtake the tagged customer, and the same holds for the customers that arrive during their service time, and so on. Finally, B denotes the tagged customer's own service requirement. If there is no customer in service upon arrival, the sojourn-time is just equal to the customer's own service requirement B . Note that the probability of arriving to a non-empty system is ρ . We thus have for S_{LNP} , the unconditional sojourn-time of the tagged customer,

$$\mathbf{P}\{S_{LNP} \leq t\} = (1 - \rho)\mathbf{P}\{B \leq t\} + \rho\mathbf{P}\{S_{LNP}(B^r) \leq t\}, \quad t \geq 0. \tag{28}$$

Here B^r denotes the (unconditional) residual service requirement of the customer in service, i.e., B^r has distribution $B^r(\cdot)$, and, hence,

$$\mathbf{P}\{S_{LNP}(B^r) \leq t\} = \int_{\tau=0}^{\infty} \mathbf{P}\{S_{LNP}(\tau) \leq t\} dB^r(\tau), \quad t \geq 0.$$

We now prove **Theorem 2.4** for non-integer $\nu > 2$ (see **Remark 4.2**), by showing that if $B(\cdot) \in \mathcal{R}_{-\nu}$ and, hence, $B^r(\cdot) \in \mathcal{R}_{-\alpha}$ with $\alpha := \nu - 1$, then $S_{\text{LNP}}(\tau)$ satisfies **Assumption 4.2** and, by **Theorem 4.1**,

$$\mathbf{P} \left\{ S_{\text{LNP}}(B^r) > \frac{x}{1 - \rho} \right\} \sim \mathbf{P}\{B^r > x\}.$$

Since $1 - B(x) = o(1 - B^r(x))$, $x \rightarrow \infty$, we have, from (28),

$$\mathbf{P} \left\{ S_{\text{LNP}} > \frac{x}{1 - \rho} \right\} \sim \rho \mathbf{P}\{B^r > x\},$$

in accordance with **Theorem 2.4**.

Remark 4.2. The case $1 < \nu < 2$ needs special treatment. As we will see below, the approach aims at verification of Condition (b) of **Assumption 4.2**, taking κ equal to the nearest integer larger than $\nu - 1$, which in this case would be $\kappa = 1$. However, for $\kappa < 2$ we cannot use (24). It is possible to verify Condition (b) using different probabilistic arguments which, however, we will not pursue here.

A different problem occurs when ν is integer-valued. The approach would aim at choosing $\kappa = \nu$. This requires that $\mathbf{E}\{B^\nu\} < \infty$, which may not be the case. (Recall that in the analysis of the $M/G/1$ PS queue we could choose κ equal to any integer larger than ν .)

In order to verify the conditions in **Assumption 4.2**, we first derive differential equations for the moments of $S_{\text{LNP}}(\tau)$. The random variable P will have the distribution of the busy-period in the $M/G/1$ queue. We assume that $m < \nu < m + 1$, for some integer $m \geq 2$, and therefore $\mathbf{E}\{B^m\} < \infty$ and $\mathbf{E}\{P^m\} < \infty$.

Remark 4.3. The fact that $\mathbf{E}\{P^m\} < \infty$ if and only if $\mathbf{E}\{B^m\} < \infty$ is a consequence of **Theorem 2.3** (since $S_{\text{LNP}} \stackrel{d}{=} P$). Note that this result was proven in [43] using Laplace-Transform techniques (cf. (14)). The same can be achieved using (probabilistic) sample-path arguments [61].

Conditioning on whether or not an arrival occurs during a time interval of length $\Delta > 0$, we obtain, for $k = 1, 2, \dots, m$,

$$\begin{aligned} & \mathbf{E}\{[S_{\text{LNP}}(\tau + \Delta)]^k\} \\ &= (1 - \lambda\Delta)\mathbf{E}\{[(\Delta + S_{\text{LNP}}(\tau))^k]\} + \lambda\Delta\mathbf{E}\{[(\Delta + S_{\text{LNP}}(\tau) + P)^k]\} + o(\Delta), \quad \Delta \rightarrow 0. \end{aligned}$$

Re-arranging terms, dividing by Δ , passing $\Delta \rightarrow 0$ and using $\mathbf{E}\{P\} = \beta/(1 - \rho)$, we obtain

$$\frac{d}{d\tau} \mathbf{E}\{S_{\text{LNP}}(\tau)^k\} = \frac{k}{1 - \rho} \mathbf{E}\{S_{\text{LNP}}(\tau)^{k-1}\} + \lambda \sum_{j=0}^{k-2} \binom{k}{j} \mathbf{E}\{S_{\text{LNP}}(\tau)^j\} \mathbf{E}\{P^{k-j}\}, \tag{29}$$

with initial condition $\mathbf{E}\{S_{\text{LNP}}(0)^k\} = \mathbf{E}\{B^k\}$. By solving (29) for $k = 1$ (setting the empty sum equal to 0), it can readily be verified that

$$\mathbf{E}\{S_{\text{LNP}}(\tau)\} = \beta + \frac{\tau}{1 - \rho},$$

which shows Condition (a) of Assumption 4.2 is satisfied. Recursively solving (29) for $k = 2, 3, \dots, m$, leads to the general form

$$\mathbf{E}\{S_{\text{LNP}}(\tau)^k\} = \left(\frac{\tau}{1-\rho}\right)^k + p_{k-1}(\tau), \tag{30}$$

where $p_{k-1}(\tau)$ denotes a polynomial in τ of degree $k - 1$. The coefficients of this polynomial can be obtained recursively by substitution into (29); in particular, $p_{k-1}(0) = \mathbf{E}\{B^k\}$. Taking $\kappa = m$ we may use (24) and (30) to show that Condition (b) of Assumption 4.2 is satisfied (for any $0 < \delta < 1$). Finally, Condition (c) can again be verified by a sample-path argument similar to that in the analysis of the $M/G/1$ PS queue.

(iii) *The M/G/1 FBPS queue.* With the FBPS discipline, the customers which so far have received the least amount of service share equally in the total capacity. Using Theorem 4.1, we will prove that the sojourn-time tail is just as heavy as the service requirement tail if $B(\cdot) \in \mathcal{R}_{-\nu}$ with $1 < \nu < 2$. (For $\nu \geq 2$ we need to study higher moments of the conditional sojourn-time.) $S_{\text{FB}}(\tau)$ denotes the sojourn-time of a customer with service requirement τ . With a straightforward sample-path argument it can be shown that $S_{\text{FB}}(\tau)$ is stochastically non-decreasing in τ .

Assuming $B(\cdot)$ is absolutely continuous, the mean and variance of the sojourn-time are given by

$$\begin{aligned} \mathbf{E}\{S_{\text{FB}}(\tau)\} &= \frac{\tau}{1-\lambda h_1(\tau)} + \frac{\lambda h_2(\tau)}{2(1-\lambda h_1(\tau))^2}, \\ \text{Var}\{S_{\text{FB}}(\tau)\} &= \frac{\lambda h_3(\tau)}{3(1-\lambda h_1(\tau))^3} + \frac{\lambda \tau h_2(\tau)}{(1-\lambda h_1(\tau))^3} + \frac{3(\lambda h_2(\tau))^2}{4(1-\lambda h_1(\tau))^4}, \end{aligned}$$

cf. [58, Form. (6.2) and (6.3)]. The functions $h_j(\tau)$, $j = 1, 2, 3$, are given by

$$h_j(\tau) = j \int_{x=0}^{\tau} x^{j-1} (1 - B(x)) \, dx.$$

These expressions can again be used [47] to prove that, for all $\varepsilon > 0$,

$$\mathbf{E}\{S_{\text{FB}}(\tau)\} \sim \frac{\tau}{1-\rho}, \quad \text{Var}\{S_{\text{FB}}(\tau)\} = o(\tau^{3-\nu+\varepsilon}), \quad \tau \rightarrow \infty.$$

Consequently, Assumption 4.2 is implied by Assumption 4.1 (choosing $\kappa = 2$ and $0 < \delta < \nu - 1$) and we may again apply Theorem 4.1 to show that $\mathbf{P}\{S_{\text{FB}} > x/(1-\rho)\} \sim \mathbf{P}\{B > x\}$, cf. Theorem 2.5.

(iv) *The M/G/1 SRPTF queue.* Now we consider an $M/G/1$ queue in which the total service capacity is always allocated to the customer with the shortest remaining processing time (Shortest-Remaining-Processing-Time-First). The service of a customer is preempted when a new customer arrives with a service requirement smaller than the remaining service requirement of the customer being served. The service of the customer that is preempted is resumed as soon as there are no other customers with a smaller amount of work in the system.

As in the $M/G/1$ FBPS queue, we restrict ourselves to the case $B(\cdot) \in \mathcal{R}_{-\nu}$ with $1 < \nu < 2$. We further assume that $B(\cdot)$ is a continuous function, hence, with probability 1, no two customers in the system have the same remaining service requirement, see [54]. The sojourn-time can be decomposed into two different periods: the waiting-time (the time until the customer is first taken into service) and the residence time (the remainder of the sojourn-time). The residence time may contain service preemption

periods caused by customers with a smaller service requirement. For a customer with service requirement τ , we denote the waiting-time by $W(\tau)$ and the residence time by $R(\tau)$. Thus, the sojourn-time is given by $S_{SR}(\tau) = W(\tau) + R(\tau)$. We define $\rho(\tau)$ as the traffic load of customers with an amount of work less than or equal to τ ,

$$\rho(\tau) := \lambda \int_{t=0}^{\tau} t \, dB(t).$$

The first two moments of $W(\tau)$ are given by

$$\begin{aligned} \mathbf{E}\{W(\tau)\} &= \lambda \frac{\int_{t=0}^{\tau} t^2 \, dB(t) + \tau^2(1 - B(\tau))}{2(1 - \rho(\tau))^2}, \\ \mathbf{E}\{W(\tau)^2\} &= \lambda \frac{\int_{t=0}^{\tau} t^3 \, dB(t) + \tau^3(1 - B(\tau))}{3(1 - \rho(\tau))^3} + \lambda^2 \int_{t=0}^{\tau} t^2 \, dB(t) \frac{\int_{t=0}^{\tau} t^2 \, dB(t) + \tau^2(1 - B(\tau))}{(1 - \rho(\tau))^4}, \end{aligned}$$

and the mean and variance of $R(\tau)$ by

$$\mathbf{E}\{R(\tau)\} = \int_{t=0}^{\tau} \frac{1}{1 - \rho(t)} dt, \quad \text{Var}\{R(\tau)\} = \lambda \int_{t=0}^{\tau} \frac{\int_{u=0}^t u^2 \, dB(u)}{(1 - \rho(t))^3} dt,$$

cf. [54]. These expressions may again be used [47] to verify that, when $1 < \nu < 2$, Assumption 4.1 implies Assumption 4.2 and, hence, $\mathbf{P}\{S_{SR} > x/(1 - \rho)\} \sim \mathbf{P}\{B > x\}$, cf. Theorem 2.6.

5. Sample-path techniques

In the present section we describe how sample-path techniques may be used to determine the tail asymptotics of the delay distribution in the $M/G/1$ queue for various disciplines. By definition, the tail distribution of a random variable reflects the occurrence of rare events. Large-deviations theory suggests that, given that a rare event occurs, it happens with overwhelming probability in the most likely way. In case light-tailed processes are involved, the most likely path typically consists of an extremely long sequence of slightly unusual events, which conspire to make the rare event under consideration occur, see for instance Anantharam [3]. In contrast, for heavy-tailed characteristics, the most likely scenario usually involves just a single catastrophic event (or generally, a ‘minimal combination’ of disastrous events that is required to cause the event under consideration to happen). Typically, the scenario entails the arrival of a customer with an exceedingly large service requirement.

The fact that the most likely scenario usually involves just a single exceptional event, provides a heuristic method for obtaining the tail asymptotics by simply computing the probability of that scenario occurring. By way of illustration, we now sketch a heuristic derivation of the tail asymptotics of the workload V in the $M/G/1$ queue as described in Section 2.

Let us focus on the workload in the system at time $t = 0$. The assumption is that a large workload level is most likely due to the prior arrival of a customer with a large service requirement B , let us say at time $t = -y$. (Of course, this assumption is nothing but an educated guess at this stage. However, it turns out that this supposition leads to the correct result, and can actually be strengthened into a rigorous proof, as will be illustrated below.) Note that from time $t = -y$ onward, the workload decreases in a roughly linear fashion at rate $1 - \rho$. So in order for the workload at time $t = 0$ to exceed the level x , the

service requirement B must be larger than $x + y(1 - \rho)$. Observing that customers arrive according to a Poisson process of rate λ , integrating w.r.t. y , and making the substitution $z = x + y(1 - \rho)$, we obtain, for large x ,

$$\mathbf{P}\{V > x\} \approx \int_{y=0}^{\infty} \mathbf{P}\{B > x + y(1 - \rho)\} \lambda \, dy = \frac{\lambda}{1 - \rho} \int_{z=x}^{\infty} \mathbf{P}\{B > z\} \, dz = \frac{\rho}{1 - \rho} \mathbf{P}\{B^r > x\}. \quad (31)$$

With some additional effort, the heuristic derivation can often be made rigorous. The typical approach consists of deriving lower and upper bounds which asymptotically coincide. It is often relatively straightforward to convert the heuristic arguments into a strict lower bound, by simply calculating the probability of the most likely scenario occurring. The construction of a suitable upper bound tends to be more challenging. The upper bound usually contains a dominant term, which corresponds to the probability of the most likely scenario. The main difficulty lies in showing that this scenario is indeed the only plausible one, in the sense that all other possible sample paths do not significantly contribute. This is done by grouping all other sample paths into a few events which must then all be shown to have an asymptotically negligible probability.

Although the above approach is fairly typical, it is hard to describe a universal method that can be mechanically executed. The identification of the most likely scenario requires some sort of an educated guess. Besides, categorizing the ‘irrelevant’ sample paths is problem-specific and far from automatic. The next lemma however characterizes the structure that typically emerges.

Lemma 5.1. *Suppose that for any $\delta > 0, \epsilon > 0, M > 0$,*

$$\mathbf{P}\{X > x\} \geq F(-\delta) \mathbf{P}\{Y > G(\epsilon)x\} \prod_{i=1}^K \mathbf{P}\{D_i^{-\delta, \epsilon}(x)\}, \quad (32)$$

$$\mathbf{P}\{X > x\} \leq F(\delta) \mathbf{P}\{Y > G(-\epsilon)x\} + \sum_{j=1}^L \mathbf{P}\{E_j^{\delta, -\epsilon}(M, x)\}, \quad (33)$$

$\mathbf{P}\{Y > x\}$ is regularly varying of index $-\nu$, $\lim_{\delta \rightarrow 0} F(\delta) = F$, $\lim_{\epsilon \rightarrow 0} G(\epsilon) = G$, $\mathbf{P}\{D_i^{-\delta, \epsilon}(x)\} \rightarrow 1$ as $x \rightarrow \infty$, and

$$\lim_{M \rightarrow \infty} \limsup_{x \rightarrow \infty} \frac{\mathbf{P}\{E_j^{\delta, -\epsilon}(M, x)\}}{\mathbf{P}\{Y > Gx\}} = 0. \quad (34)$$

Then

$$\mathbf{P}\{X > x\} \sim F \mathbf{P}\{Y > Gx\}.$$

Proof. The proof is straightforward. Relying on the lower bound (32) and the fact that $\mathbf{P}\{D_i^{-\delta, \epsilon}(x)\} \rightarrow 1$ as $x \rightarrow \infty$, we obtain

$$\liminf_{x \rightarrow \infty} \frac{\mathbf{P}\{X > x\}}{F \mathbf{P}\{Y > Gx\}} \geq \frac{F(-\delta)}{F} \liminf_{x \rightarrow \infty} \frac{\mathbf{P}\{Y > G(\epsilon)x\}}{\mathbf{P}\{Y > Gx\}}.$$

Letting $\delta, \epsilon \downarrow 0$, and recalling that $\mathbf{P}\{Y > x\}$ is regularly varying, we find

$$\liminf_{x \rightarrow \infty} \frac{\mathbf{P}\{X > x\}}{F \mathbf{P}\{Y > Gx\}} \geq 1.$$

Similarly, using the upper bound (33) and (34), and observing that $\mathbf{P}\{Y > x\}$ is regularly varying of index $-\nu$, we deduce

$$\limsup_{x \rightarrow \infty} \frac{\mathbf{P}\{X > x\}}{F\mathbf{P}\{Y > Gx\}} \leq \frac{F(\delta)}{F} \limsup_{x \rightarrow \infty} \frac{\mathbf{P}\{Y > G(-\epsilon)x\}}{\mathbf{P}\{Y > Gx\}}.$$

Letting $\delta, \epsilon \downarrow 0$, we conclude

$$\limsup_{x \rightarrow \infty} \frac{\mathbf{P}\{X > x\}}{F\mathbf{P}\{Y > Gx\}} \leq 1. \quad \square$$

It is worth observing that, like in Section 4, the above proof technique in fact allows for intermediately regularly varying distributions.

The events $E_j^{\delta, -\epsilon}(M, x)$ may be interpreted as events leading to $\{X > x\}$, other than the ‘typical’ event. Often, $E_j^{\delta, -\epsilon}(M, x)$ is independent of M and has the simpler property $\mathbf{P}\{E_j^{\delta, -\epsilon}\} = o(\mathbf{P}\{Y > Gx\})$ as $x \rightarrow \infty$. Sometimes, M is required as an additional auxiliary parameter.

As a ‘toy example’, we now sketch how the above lemma may be used to strengthen the heuristic derivation of (31) into a rigorous proof. The approach is similar as outlined in Chapter 2 of Zwart [60]. We use the time-reversed sample-path representation

$$V \stackrel{d}{=} \sup_{t \geq 0} \{A(0, t) - t\} \tag{35}$$

with $A(0, t)$ denoting the amount of work arriving in the time interval $(0, t)$.

Our ‘educated guess’ is that (31) provides, indeed, the correct asymptotics.

We first construct a lower bound of the form (32). For any $c < \rho$, define $U^c := \sup_{t \geq 0} \{ct - A(0, t)\}$. For any $\delta > 0, \epsilon > 0$,

$$\begin{aligned} \mathbf{P}\{V > x\} &\geq \int_{y=0}^{\infty} \mathbf{P}\{A(0, y) + B - y > x\} \lambda \, dy \\ &\geq \lambda \int_{y=0}^{\infty} \mathbf{P}\{A(0, y) - y(\rho - \delta) \geq -\epsilon x\} \mathbf{P}\{B > x(1 + \epsilon) + y(1 - \rho + \delta)\} \, dy \\ &\geq \mathbf{P}\left\{\inf_{u \geq 0} \{A(0, u) - u(\rho - \delta)\} \geq -\epsilon x\right\} \lambda \int_{y=0}^{\infty} \mathbf{P}\{B > x(1 + \epsilon) + y(1 - \rho + \delta)\} \, dy \\ &= \frac{\rho}{1 - \rho + \delta} \mathbf{P}\{B^r > x(1 + \epsilon)\} \mathbf{P}\{U^{\rho - \delta} \leq \epsilon x\}. \end{aligned}$$

Note that $\mathbf{P}\{U^{\rho - \delta} \leq \epsilon x\} \rightarrow 1$ as $x \rightarrow \infty$ because of the law of large numbers.

We now proceed to derive an upper bound of the form (33). For any interval $I \subseteq \mathbb{R}^+$, define $V(I) := \sup_{t \in I} \{A(0, t) - t\}$. For any $y \geq 0$, let $N_y(I)$ be the number of customers arriving during the time interval I whose service requirement exceeds the value y . Then, for all $M \geq 0$,

$$\begin{aligned} \mathbf{P}\{V > x\} &\leq \mathbf{P}\{V([0, Mx]) > x\} + \mathbf{P}\{V((Mx, \infty)) > x\} \\ &= \mathbf{P}\{V([0, Mx]) > x; N_{\epsilon x}([0, Mx]) = 0\} + \mathbf{P}\{V([0, Mx]) > x; N_{\epsilon x}([0, Mx]) = 1\} \\ &\quad + \mathbf{P}\{V([0, Mx]) > x; N_{\epsilon x}([0, Mx]) \geq 2\} + \mathbf{P}\{V((Mx, \infty)) > x\}. \end{aligned}$$

The second term corresponds to the only plausible scenario and is dominant. Let $\tau(\epsilon, x)$ be the arrival time of the large customer. As indicated in Section 2.4 of [60] (see also [62]), it may be shown that for any $\delta > 0, \epsilon > 0$, as $x \rightarrow \infty$,

$$\begin{aligned} & \mathbf{P}\{V([0, Mx]) > x; N_{\epsilon x}([0, Mx]) = 1\} \\ & \leq \mathbf{P}\{A(0, \tau(\epsilon, x)^-) \leq (\rho + \delta)\tau(\epsilon, x); A(0, \tau(\epsilon, x)) - \tau(\epsilon, x) \geq (1 - \delta)x\} + o(\mathbf{P}\{B^r > x\}) \\ & \leq \int_{y=0}^{\infty} \mathbf{P}\{B > x(1 - \epsilon) + y(1 - \rho - 2\delta)\} \lambda \, dy + o(\mathbf{P}\{B^r > x\}) \\ & = \frac{\rho}{1 - \rho - 2\delta} \mathbf{P}\{B^r > x(1 - \epsilon)\} + o(\mathbf{P}\{B^r > x\}). \end{aligned}$$

Applying Lemma 5.1 completes the proof, once we have shown that each of the other three terms can asymptotically be neglected.

For the first term, one may exploit a powerful lemma of Resnick and Samorodnitsky [52] to show that for any $\mu > 0$ there exists an $\epsilon > 0$ such that

$$\mathbf{P}\{V([0, Mx]) > x; N_{\epsilon x}([0, Mx]) = 0\} = o(x^{-\mu})$$

as $x \rightarrow \infty$. The idea is that when there are no large service requirements, the process $\{A(0, t) - t\}$ cannot significantly deviate from its normal drift over long intervals of the order x , so that the workload cannot reach a large level.

In order to control the third term, using the Poisson structure of the arrival process, it can be shown that $\mathbf{P}\{N_{\epsilon x}([0, Mx]) \geq 2\} = o(\mathbf{P}\{B^r > x\})$ as $x \rightarrow \infty$. In words, this means that the probability of two large service requirements occurring in a time interval of order x is asymptotically negligible compared to that of just one large service requirement.

Finally, for the fourth term, we use the upper bound (for some $\delta > 0$)

$$\mathbf{P}\{V((Mx, \infty)) > x\} \leq 2\mathbf{P}\{V^{1-2\delta} > \delta Mx\}, \tag{36}$$

with $V^{1-2\delta}$ the steady-state workload in an $M/G/1$ queue with a server working at speed $1 - 2\delta$, see p. 197 of [60] for a similar statement. The right-hand side in (36) can be upper bounded by using a result (obtained from first principles) of Mikosch [44]:

$$\mathbf{P}\{V^{1-2\delta} > x\} \leq (C_\delta + o(1))\mathbf{P}\{B^r > x\}. \tag{37}$$

This implies, using the previously derived lower bound and the upper bounds (36) and (37),

$$\lim_{M \rightarrow \infty} \limsup_{x \rightarrow \infty} \frac{\mathbf{P}\{V((Mx, \infty)) > x\}}{\mathbf{P}\{V > x\}} = 0.$$

This result indicates that overflow of a level of order x must occur ‘in linear time’, since otherwise the process $\{A(0, t) - t\}$ must deviate from its normal drift for a prohibitively long period of time.

The above proof exploits and confirms the large deviations notion that a large workload level is typically due to a single large service requirement by implicitly characterizing the most likely sample path. In the literature, similar statements have been proven by Asmussen and Klüppelberg [6], and Baccelli and Foss [9]. We note that the results in these papers rely on the fact that the workload asymptotics were already available, unlike the proof given here which provides both the asymptotics and the ‘single big jump’ result together.

Nevertheless, we emphasize that we used the above proof technique for illustration purposes only. The machinery is unnecessarily heavy for determining the workload asymptotics in the ordinary $M/G/1$ queue, for which simpler methods are available (see Sections 3 and 6). The true merits of the methodology become manifest in more complicated systems, such as fluid queues or GPS models, where typically no useful expression for the LST is available [14,19,62].

5.1. The single-class case

We now turn the attention to the tail asymptotics of the delay distribution in the $M/G/1$ queue. In contrast to the workload distribution, the delay distribution *does* strongly depend on the service discipline that is used.

(i) *The $M/G/1$ FCFS queue.* For FCFS, the waiting-time is simply equal to the workload at the time of arrival. Because of the PASTA property, it then follows from (31) that:

$$\mathbf{P}\{W_{\text{FCFS}} > x\} \sim \frac{\rho}{1-\rho} \mathbf{P}\{B^r > x\},$$

which agrees with Theorem 2.1.

(ii) *The $M/G/1$ LCFS-NP queue.* For LCFS Non-Preemptive priority, the waiting-time is equal to 0 with probability $1 - \rho$, and with probability ρ it is equal to a busy-period starting with a residual service requirement, which gives

$$\mathbf{P}\{W_{\text{LNP}} > x\} \sim \rho \mathbf{P}\{B^r > x(1 - \rho)\}, \quad (38)$$

as asserted in Theorem 2.4.

A heuristic derivation of the above formula proceeds as follows. Consider a tagged customer arriving at time $t = 0$. The assumption is that a long waiting-time is most likely due to a large service requirement B of the customer in service, if any. The waiting-time W of the tagged customer then consists of the remaining service requirement, $B - y$, plus the amount of work arriving during its own waiting-time, which is approximately ρW , so that $W \approx B - y + \rho W$, or equivalently, $W \approx (B - y)/(1 - \rho)$. So in order for the waiting-time to exceed the value x , the service requirement B must be larger than $y + x(1 - \rho)$. Thus, observing that arrivals occur as a Poisson process of rate λ , and integrating w.r.t. y , we find, for large x ,

$$\mathbf{P}\{W_{\text{LNP}} > x\} \approx \int_{y=0}^{\infty} \mathbf{P}\{B > x(1 - \rho) + y\} \lambda \, dy = \rho \mathbf{P}\{B^r > x(1 - \rho)\},$$

which is in agreement with (38). The above heuristic derivation may be translated into a rigorous proof in a similar fashion as indicated for the workload asymptotics.

(iii) *The $M/G/1$ LCFS-PR queue.* For LCFS Preemptive-Resume, the sojourn-time is simply equal to the busy-period, yielding

$$\mathbf{P}\{S_{\text{LPR}} > x\} \sim \frac{1}{1-\rho} \mathbf{P}\{B > x(1 - \rho)\},$$

as stated in Theorem 2.3. A sample-path proof of the tail asymptotics of the busy-period distribution may be found in [61].

(iv) *The M/G/1 PS queue.* We now turn to the tail asymptotics of the sojourn-time for the Processor-Sharing discipline. Consider a tagged customer arriving at time $t = 0$. The sojourn-time S of the tagged customer consists of its own service requirement B plus the amount of service provided to other customers during its sojourn-time. In case of a long sojourn-time, the amount of service received by other customers will be approximately ρS , so that $S \approx B + \rho S$, or equivalently, $S \approx B/(1 - \rho)$. The assumption is thus that a long sojourn-time is most likely due to a large service requirement of the tagged customer itself, suggesting that, for large x ,

$$\mathbf{P}\{S_{PS} > x\} \sim \mathbf{P}\{B > x(1 - \rho)\}, \tag{39}$$

which corroborates with [Theorem 2.2](#).

We now show how the above rough derivation may be used as the basis for a rigorous proof of (39) using lower and upper bounds along the lines of [Lemma 5.1](#). The proof is similar to that in [34]. Let B_0 and S_0 be the service requirement and the sojourn-time, respectively, of a tagged customer arriving at time $t = 0$. Let B_i and T_i denote the service requirement and the arrival time of the i th customer arriving after time $t = 0$. Let $L(0)$ be the number of customers in the system just before time $t = 0$, and let B_l^r denote the remaining service requirement of the l th customer. We use the sample-path representation

$$S_0 = B_0 + \sum_{l=1}^{L(0)} \min\{B_l^r, B_0\} + \sum_{i=1}^{N((0, S_0))} \min\{B_i, R_0(T_i)\}, \tag{40}$$

with $N((0, t))$ denoting the number of customers arriving during the time interval $(0, t)$, and with $R_0(t)$ representing the remaining service requirement of the tagged customer at time t .

The next lemma presents a lower bound for the sojourn-time of the tagged customer. Denote $Z(t) := \sum_{i=1}^{N((0, t))} \max\{B_i - R_0(T_i), 0\}$.

Lemma 5.2. *For any $\delta > 0$,*

$$S_0(1 - \rho + \delta) \geq B_0 - U^{\rho-\delta} - Z(S_0).$$

Proof. Using the representation (40), we have

$$\begin{aligned} S_0(1 - \rho + \delta) &= B_0 + \sum_{l=1}^{L(0)} \min\{B_l^r, B_0\} + \sum_{i=1}^{N((0, S_0))} \min\{B_i, R_0(T_i)\} - (\rho - \delta)S_0 \\ &\geq B_0 + \sum_{i=1}^{N((0, S_0))} B_i - (\rho - \delta)S_0 + \sum_{i=1}^{N((0, S_0))} \min\{B_i, R_0(T_i)\} - \sum_{i=1}^{N((0, S_0))} B_i \\ &= B_0 + A(0, S_0) - (\rho - \delta)S_0 + \sum_{i=1}^{N((0, S_0))} \min\{R_0(T_i) - B_i, 0\} \\ &\geq B_0 + \inf_{t \geq 0} \{A(0, t) - (\rho - \delta)t\} - \sum_{i=1}^{N((0, S_0))} \max\{B_i - R_0(T_i), 0\} \\ &= B_0 - U^{\rho-\delta} - Z(S_0). \end{aligned} \quad \square$$

The next lemma provides an upper bound for the sojourn-time of the tagged customer. For any $y > 0$, let $A_y(0, t)$ be a version of the process $A(0, t)$ where all service requirements of arriving customers are truncated at the level y . For any $c > \rho$, define $V_y^c := \sup_{t \geq 0} \{A_y(0, t) - ct\}$.

Lemma 5.3. For any $\delta > 0$,

$$(1 - \rho - \delta)S_0 \leq \hat{B}_0 + V_{B_0}^{\rho+\delta},$$

with $\hat{B}_0 := B_0 + \sum_{l=1}^{L(0)} \min\{B_l^r, B_0\}$.

Proof. Using the representation (40),

$$\begin{aligned} S_0(1 - \rho - \delta) &= B_0 + \sum_{l=1}^{L(0)} \min\{B_l^r, B_0\} + \sum_{i=1}^{N((0, S_0))} \min\{B_i, R_0(T_i)\} - (\rho + \delta)S_0 \\ &\leq \hat{B}_0 + \sum_{i=1}^{N((0, S_0))} \min\{B_i, B_0\} - (\rho + \delta)S_0 \\ &= \hat{B}_0 + A_{B_0}(0, S_0) - (\rho + \delta)S_0 \leq \hat{B}_0 + \sup_{t \geq 0} \{A_{B_0}(0, t) - (\rho + \delta)t\} \\ &= \hat{B}_0 + V_{B_0}^{\rho+\delta}. \end{aligned} \quad \square$$

The above two lemmas provide the necessary ingredients for the proof of (39) along the lines of Lemma 5.1.

Proof of Theorem 2.2. Lower bound. Using Lemma 5.2, noting that $S_0 \geq B_0$, we obtain

$$\begin{aligned} \mathbf{P}\{S_{PS} > x\} &\geq \mathbf{P}\{B_0 - U^{\rho-\delta} - Z(S_0) > (1 - \rho + \delta)x\} \\ &\geq \mathbf{P}\{B_0 > (1 - \rho + \delta + 2\epsilon)x\} \mathbf{P}\{U^{\rho-\delta} \leq \epsilon x\} \inf_{y \geq (1-\rho+\delta+2\epsilon)x} \mathbf{P}\{Z(y) > \epsilon x\}. \end{aligned}$$

Because of the law of large numbers, $\mathbf{P}\{U^{\rho-\delta} \leq \epsilon x\} \rightarrow 1$ as $x \rightarrow \infty$. As observed in [34], $\inf_{y \geq (1-\rho+\delta+2\epsilon)x} \mathbf{P}\{Z(y) > \epsilon x\} \rightarrow 1$ as $x \rightarrow \infty$.

Upper bound. Using Lemma 5.3, we find

$$\mathbf{P}\{S_{PS} > x\} \leq \mathbf{P}\{\hat{B}_0 + V_{B_0}^{\rho+\delta} > (1 - \rho - \delta)x\} \leq \mathbf{P}\{\hat{B}_0 > (1 - \rho - \delta - \epsilon)x\} + \mathbf{P}\{V_{B_0}^{\rho+\delta} > \epsilon x\}.$$

As demonstrated in [34], $\mathbf{P}\{V_{B_0}^{\rho+\delta} > \epsilon x\} = o(\mathbf{P}\{B > x\})$ as $x \rightarrow \infty$, and $\mathbf{P}\{\hat{B}_0 > x\} \sim \mathbf{P}\{B > x\}$.

Invoking Lemma 5.1 then completes the proof. □

5.2. The multi-class case

We now consider the tail asymptotics of the waiting-time in the multi-class $M/G/1$ queue with priorities as described in Section 3.2. We focus on the case of a Non-Preemptive priority discipline. As mentioned in Section 3.2, the tail asymptotics in the case of a Preemptive-Resume policy immediately follow from the results for a high-priority class in isolation and a low-priority class in the Non-Preemptive priority

scenario. We assume that the service requirement distribution of at least one of the classes has a regularly varying tail. Let \mathcal{M} be the index set of the classes with the ‘heaviest’ tail.

Consider a tagged class- k customer arriving at time $t = 0$. The assumption is that a long waiting-time is typically due to the prior arrival of a customer with a large service requirement B , let us say at time $t = -y$, which may belong to any of the classes $m \in \mathcal{M}$. Of course, how likely it is for the culprit customer to belong to a given class $m \in \mathcal{M}$ depends on the arrival rates and mean service requirements of the various classes. Due to the Non-Preemptive priority policy, the identity of the culprit customer is not of any relevance for the impact on the tagged customer. However, the effect does strongly depend on the identity of the tagged customer itself. For compactness, denote $\sigma_k = \sum_{l=1}^k \rho_l$. Note that from time $t = -y$ onward, the amount of work in the system that has precedence over the service of the tagged customer decreases in a roughly linear fashion at rate $1 - \sigma_k$. In addition, the tagged customer must wait for the amount of work arriving during its own waiting-time W_k from higher-priority classes at rate σ_{k-1} . Thus, $W_k \approx B - y(1 - \sigma_k) + W_k \sigma_{k-1}$. So in order for the waiting-time of the tagged customer to exceed the value x , the service requirement B must be larger than $x(1 - \sigma_{k-1}) + y(1 - \sigma_k)$. Observing that class- m customers arrive as a Poisson process of rate λ_m , integrating w.r.t. y , and making the substitution $z = x(1 - \sigma_{k-1}) + y(1 - \sigma_k)$, we obtain, for large x ,

$$\begin{aligned} \mathbf{P}\{W_k > x\} &\approx \sum_{m \in \mathcal{M}} \int_{y=0}^{\infty} \mathbf{P}\{B_m > x(1 - \sigma_{k-1}) + y(1 - \sigma_k)\} \lambda_m \, dy \\ &= \sum_{m \in \mathcal{M}} \frac{\rho_m}{1 - \sigma_k} \mathbf{P}\{B'_m > x(1 - \sigma_{k-1})\}. \end{aligned} \tag{41}$$

6. Subexponential asymptotics and random sums

In this section, we relax the assumption of a regularly varying distribution function, and focus on the more general case of the $M/G/1$ queue with a *subexponential* service requirement distribution. The class of subexponential distributions is defined as follows.

Definition 6.1. A distribution function $F(x) = \mathbf{P}\{X \leq x\}$ is subexponential ($F(\cdot) \in \mathcal{S}$) if

$$\mathbf{P}\{X_1 + \dots + X_n > x\} \sim n\mathbf{P}\{X_1 > x\},$$

for any $n \geq 2$, with X_1, \dots, X_n i.i.d. copies of X .

Sometimes, the random variable X (rather than its distribution function $F(\cdot)$) is called subexponential. Subexponential distributions have been introduced by Chistyakov [28]. Note that subexponentiality is equivalent to

$$\mathbf{P}\{X_1 + \dots + X_n > x\} \sim \mathbf{P}\{\max_{i=1, \dots, n} X_i > x\}.$$

Thus, large values of sums of subexponential random variables have the appealing property that they are dominated by their largest term.

As mentioned earlier, the class of subexponential distributions is larger than the class of regularly varying distributions. Examples of subexponential distributions which are not regularly varying, are the lognormal distribution and Weibull distributions with tails of the form e^{-x^β} , $0 < \beta < 1$.

In extending results from the regularly varying case to the subexponential case, one faces several difficulties. First of all, there is no characterization of subexponential distribution functions in terms of their LST, which rules out the techniques of Section 3. A further complication is that, sometimes, the class seems too large to work with. Often, one has to invoke additional regularity conditions; see the examples below.

An exception though, is when an explicit expression for the distribution function is available, in particular when the random variable of interest can be expressed as a *random sum*. If a random-sum representation is not available, one has to resort to refinements of the method given in Section 5; specific references are given below.

6.1. The single-class case

(i) *The M/G/1 FCFS queue.* The geometric structure of the LST given by (7) allows for an explicit inversion, leading to

$$\mathbf{P}\{W > x\} = (1 - \rho) \sum_{n=0}^{\infty} \rho^n \mathbf{P}\{B_1^r + \dots + B_n^r > x\}, \tag{42}$$

which can equivalently be phrased as

$$W \stackrel{d}{=} B_1^r + \dots + B_N^r, \tag{43}$$

with N a geometrically distributed random variable with parameter ρ . Therefore, (42) is called a (geometric) random-sum representation.

From the definition, it is clear that subexponential random variables are well-suited to analyze random sums. If one assumes that B^r is subexponential, one obtains

$$\begin{aligned} \mathbf{P}\{W > x\} &= (1 - \rho) \sum_{n=0}^{\infty} \rho^n \mathbf{P}\{B_1^r + \dots + B_n^r > x\} \sim (1 - \rho) \sum_{n=0}^{\infty} \rho^n n \mathbf{P}\{B^r > x\} \\ &= \frac{\rho}{1 - \rho} \mathbf{P}\{B^r > x\}. \end{aligned}$$

This derivation assumes that interchanging limit and summation is allowed, which is guaranteed by the following upper bound, due to Kesten (see [7]): if B^r is subexponential and $\epsilon > 0$, then there exists a constant $K = K_\epsilon$ such that

$$\mathbf{P}\{B_1^r + \dots + B_n^r > x\} \leq K(1 + \epsilon)^n \mathbf{P}\{B_1^r > x\}.$$

Now, use this bound with ϵ sufficiently small that $\rho(1 + \epsilon) < 1$. The validity of the above procedure then follows from the dominated convergence theorem.

Note that this derivation assumes that B^r (rather than B itself) is subexponential. The question whether subexponentiality of B implies that of B^r is an open problem, but the implication can be shown to hold in all cases of practical interest. For example, if B is regularly varying with index $-\nu$, then B^r is regularly varying with index $1 - \nu$, as shown in Section 3.

The above result $\mathbf{P}\{W > x\} \sim \rho/(1 - \rho) \mathbf{P}\{B^r > x\}$ was first obtained by Pakes [49] for the more general GI/G/1 queue (in this more general case, one still has a random-sum representation for W). In addition, Korshunov [41] established a converse result

$$B^r \in \mathcal{S} \Leftrightarrow W \in \mathcal{S} \Leftrightarrow \mathbf{P}\{W > x\} \sim \frac{\rho}{1 - \rho} \mathbf{P}\{B^r > x\}. \tag{44}$$

This result (valid for the $GI/G/1$ queue) reveals a deep connection between subexponentiality and random sums.

(ii) *The $M/G/1$ LCFS-PR queue.* As mentioned in Section 3, the sojourn-time distribution in the case of LCFS Preemptive-Resume is exactly equal to the busy-period distribution. Unfortunately, no suitable random-sum representation is available in this case.

Nevertheless, asymptotics for the busy-period under subexponentiality have recently been obtained by Jelenković and Momčilović [35], and by Baltrunas et al. [8]. In both studies, the asymptotics

$$\mathbf{P}\{P > x\} \sim \frac{1}{1 - \rho} \mathbf{P}\{B > (1 - \rho)x\}$$

are shown to hold under some additional smoothness conditions on the service requirement distribution (besides subexponentiality). Notably, this asymptotic form *fails to hold* when the distribution of B has a Weibull tail of the form e^{-x^β} , with $1/2 \leq \beta < 1$. In particular, a necessary condition is that the tail distribution of B is *square-root insensitive* [36]:

$$\mathbf{P}\{B > x\} \sim \mathbf{P}\{B > x - \sqrt{x}\}. \tag{45}$$

An explanation of this phenomenon is given below. The proof in [35] is based on an extension of the method in Section 5. The approach in [8] exploits a Spitzer identity for first-passage times of random walks.

(iii) *The $M/G/1$ PS queue.* An extension of Theorem 2.2 to the class of subexponential service requirements has recently been established by Jelenković and Momčilović [34]. Their proof can be viewed as an extension of the methods of Section 5. Again, (45) is shown to be a necessary condition for the tail equivalence

$$\mathbf{P}\{S_{PS} > x\} \sim \mathbf{P}\{B > (1 - \rho)x\}.$$

The general idea behind the necessity of (45) is the following: the probability of the rare event to happen should be such that its asymptotic behavior is invariant for random fluctuations governed by the Central Limit Theorem (CLT). In the particular case of Processor Sharing, this can be illustrated as follows: the typical rare event $\{S_{PS} > x\}$ is determined by the event $\{B > x(1 - \rho)\}$. Define the inverse process $S_{PS}^{\leftarrow}(x)$ of $S_{PS}(x)$ by

$$S_{PS}^{\leftarrow}(x) = \inf\{t : S_{PS}(t) \geq x\}.$$

It can be shown that, due to the CLT, one has $S_{PS}^{\leftarrow}(x) = (1 - \rho)x + O(\sqrt{x})$. Hence,

$$\mathbf{P}\{S_{PS} > x\} = \mathbf{P}\{B > S_{PS}^{\leftarrow}(x)\} \approx \mathbf{P}\{B > (1 - \rho)x + O(\sqrt{x})\}, \tag{46}$$

which explains Condition (45). Note that this condition is always satisfied if B is regularly varying. On the other hand, this example shows that the heuristics described in Section 5 require caution when considering the full class of subexponential distributions.

6.2. The multi-class case

In this section, we analyze the tail behavior of the low-priority waiting-time distribution $\mathbf{P}\{W_2 > x\}$ for the priority queue with two classes. W_2 is defined to be the time from arrival until the start of the

(later possibly interrupted) service. Note that W_2 has the same distribution for both the Preemptive and Non-Preemptive case.

Abate and Whitt [2] derive the following random-sum representation for the low-priority waiting-time distribution:

$$W_2 \stackrel{d}{=} Y_1 + \dots + Y_N, \tag{47}$$

with N a geometric random variable with parameter ρ , independent of the i.i.d. sequence $Y_i, i \geq 1$, whose distribution function can be expressed as

$$\mathbf{P}\{Y_1 \leq x\} = \frac{\rho_1}{\rho} H_1(x) + \frac{\rho_2}{\rho} H_2(x).$$

As shown, in [2], the function $H_1(x)$ is determined by the residual busy-period distribution of the high-priority class. For our purposes, the following random-sum characterization is convenient:

$$\mathbf{P}\{P_1^r \leq x\} = (1 - \rho_1) \sum_{n=0}^{\infty} \rho_1^n H_1^{(n+1)*}(x). \tag{48}$$

The function $H_2(x)$ is the distribution function of a busy-period of class-1 customers, with an *exceptional first service* B_2^r , which we denote by $P_1(B_2^r)$. Such a busy-period has the following representation:

$$P_1(B_2^r) \stackrel{d}{=} B_2^r + \sum_{i=1}^{N_1(B_2^r)} P_{1,i}. \tag{49}$$

In this expression, $N_1(\cdot)$ is a Poisson process of rate λ_1 , and $P_{1,i}, i \geq 1$ are i.i.d. copies of a high-priority busy-period.

Based upon the representations (48) and (49), one can derive the tail behavior of $\mathbf{P}\{Y_1 > x\}$, which in conjunction with the random-sum representation (47), leads to the tail asymptotics of $\mathbf{P}\{W_2 > x\}$. To illustrate this, we focus on the following two special cases: (i) class 1 subexponential, class 2 light-tailed; (ii) class 2 subexponential, class 1 light-tailed. In both cases, we assume that the condition (45), as well as some technical conditions stated in [34], are satisfied.

First, assume that B_1, B_1^r are subexponential and B_2 is light-tailed. It can be shown using [35], that the residual busy-period P_1^r is then subexponential as well. In particular,

$$\mathbf{P}\{P_1^r > x\} \sim \frac{1}{1 - \rho_1} \mathbf{P}\{B_1^r > (1 - \rho_1)x\}. \tag{50}$$

This implies that P_1^r is subexponential, since subexponentiality is closed under tail-equivalence, see, e.g. [49]. But then, using the reverse implication of Pakes' theorem (44) and the random-sum representation (48), we must also have $H_1(\cdot) \in \mathcal{S}$, and

$$\mathbf{P}\{P_1^r > x\} \sim \frac{1}{1 - \rho_1} \bar{H}_1(x), \tag{51}$$

with $\bar{H}_1(x) = 1 - H_1(x)$. Combining (50) and (51) yields

$$\bar{H}_1(x) \sim \mathbf{P}\{B_1^r > (1 - \rho_1)x\}. \tag{52}$$

The tail $\bar{H}_2(x)$ can be derived from the random-sum representation (49). Even though this is not a geometric random sum, the conclusion regarding its asymptotic behavior remains the same, noting that B_2^r and $N_1(B_2^r)$ are light-tailed. We conclude that

$$\bar{H}_2(x) \sim \mathbf{E}\{N_1(B_2^r)\} \mathbf{P}\{P_{1,1} > x\}. \tag{53}$$

From this expression, it can be shown that $\bar{H}_2(x) = o(\bar{H}_1(x))$. We conclude that

$$\mathbf{P}\{Y_1 > x\} \sim \frac{\rho_1}{\rho} \bar{H}_1(x) \sim \frac{\rho_1}{\rho} \mathbf{P}\{B_1^r > (1 - \rho_1)x\}.$$

Combining this with the random-sum representation for W_2 , we conclude that, if B_1, B_1^r are subexponential and B_2 light-tailed, then

$$\mathbf{P}\{W_2 > x\} \sim \frac{\rho_1}{1 - \rho} \mathbf{P}\{B_1^r > (1 - \rho_1)x\}. \tag{54}$$

Next, we consider the opposite case where class 1 has light-tailed characteristics and $B_2^r \in \mathcal{S}$. In this case, $\bar{H}_1(x)$ can be shown to have exponential decay; the dominant term is $\bar{H}_2(x) = \mathbf{P}\{P_1(B_2^r) > x\}$. In turn, this tail is completely determined by the tail behavior of B_2^r : since $N_1(\cdot)$ and $P_{1,i}$ are both light-tailed, it can be shown that they only contribute to $P_1(B_2^r)$ through their means. This can be rigorously justified by using results in [36]. These considerations imply

$$\bar{H}_2(x) = \mathbf{P}\{P_1(B_2^r) > x\} \sim \mathbf{P}\{(1 + \lambda_1 \mathbf{E}\{P_{1,1}\})B_2^r > x\} = \mathbf{P}\{B_2^r > (1 - \rho_1)x\}.$$

Hence,

$$\mathbf{P}\{Y_1 > x\} \sim \frac{\rho_2}{\rho} \bar{H}_2(x) \sim \frac{\rho_2}{\rho} \mathbf{P}\{B_2^r > (1 - \rho_1)x\}.$$

Finally, using the random-sum representation for W_2 , we obtain

$$\mathbf{P}\{W_2 > x\} \sim \frac{\rho_2}{1 - \rho} \mathbf{P}\{B_2^r > (1 - \rho_1)x\}. \tag{55}$$

Note that both (54) and (55) agree with the expression (41) given in Section 5, since the latter expression reduces to a single term when only one of the service requirement distributions is heavy-tailed.

7. Conclusion

In this paper, we have surveyed the tail behavior of the waiting-time and/or sojourn-time distributions for several service disciplines. It turns out that, if the service time distribution is regularly varying with index $-\nu$, the waiting-time distribution in FCFS and LCFS-NP is heavier, viz. regularly varying with index $1 - \nu$. This is in contrast with service disciplines like PS, FBPS, SRPTF, and LCFS-PR: these disciplines all yield a sojourn-time tail of index $-\nu$. These results are reviewed in Section 2, and proved in several different ways in Sections 3–6, where also multiclass disciplines are treated.

The results in this paper raise several questions. First of all, one wonders whether (given a service time distribution that is regularly varying of index $-\nu$) the only possible indices of the sojourn-time distribution in a work-conserving single-server queue are $-\nu$ and $1 - \nu$. We believe that this is not the case; this is a topic for future research.

Another important question is the validity of the asymptotics for moderate values. Results in [1] for the FCFS queue show that the asymptotics may behave poorly as approximation. Moreover, asymptotic estimates tend to *underestimate* the true value of the exceedance probability. An explanation of this phenomenon is hidden in Section 5 of the present paper: the asymptotics are fully driven by the most likely scenario. However, other scenarios may be relevant as well for moderate values of x . To speed up the convergence of the asymptotic approximations to the true value of the sojourn-time tail probability, one could try to obtain more terms in the expansion. This is a problem which is largely open (an exception is FCFS, see, e.g. [1]). Another alternative to get sojourn-time tail probabilities is numerical transform inversion: transforms of the distributions are available for many service disciplines, see Section 3 of the present paper.

References

- [1] J. Abate, G.L. Choudhury, W. Whitt, Waiting-time tail probabilities in queues with long-tail service-time distributions, *Queue. Syst.* 16 (1994) 311–338.
- [2] J. Abate, W. Whitt, Asymptotics for $M/G/1$ low-priority waiting-time tail probabilities, *Queue. Syst.* 25 (1997) 173–233.
- [3] V. Anantharam, How large delays build up in a $GI/G/1$ queue, *Queue. Syst.* 5 (1988) 345–368.
- [4] V. Anantharam, Scheduling strategies and long-range dependence, *Queue. Syst.* 33 (1999) 73–89.
- [5] A. Arvidsson, P. Karlsson, On traffic models for TCP/IP, in: P. Key, D. Smith (Eds.), *Teletraffic Engineering in a Competitive World*, Proceedings of the ITC-16, Edinburgh, UK, North-Holland, Amsterdam, 1999, pp. 457–466.
- [6] S. Asmussen, C. Klüppelberg, Stationary $M/G/1$ excursions in the presence of heavy tails, *J. Appl. Probab.* 33 (1996) 208–212.
- [7] K.B. Athreya, P.E. Ney, *Branching Processes*, Springer, Berlin, 1972.
- [8] A. Baltrunas, D.J. Daley, C. Klüppelberg, Tail behaviour of the busy-period of a $GI/G/1$ queue with subexponential service times, Technical Report, Munich University of Technology, 2002.
- [9] F. Baccelli, S. Foss, Moments and tails in monotone-separable stochastic networks, Research Report RR 4197, INRIA Rocquencourt, 2001.
- [10] J. Beran, R. Sherman, M.S. Taqqu, W. Willinger, Long-range dependence in variable-bit-rate video traffic, *IEEE Trans. Commun.* 43 (1995) 1566–1579.
- [11] N.H. Bingham, R.A. Doney, Asymptotic properties of super-critical branching processes. I: The Galton–Watson process, *Adv. Appl. Probab.* 6 (1974) 711–731.
- [12] N.H. Bingham, C.M. Goldie, J.L. Teugels, *Regular Variation*, Cambridge University Press, Cambridge, UK, 1987.
- [13] S.C. Borst, O.J. Boxma, P.R. Jelenković, Coupled processors with regularly varying service times, in: Proceedings of the IEEE Infocom, Tel-Aviv, Israel, 2000, pp. 157–164.
- [14] S.C. Borst, O.J. Boxma, P.R. Jelenković, Reduced-load equivalence and induced burstiness in GPS queues with long-tailed traffic flows, *Queue. Syst.* 43 (2003) 273–306.
- [15] S.C. Borst, O.J. Boxma, J.A. Morrison, R. Núñez-Queija, The equivalence between processor sharing and service in random order, *Oper. Res. Lett.* 31 (2003) 254–262.
- [16] S.C. Borst, O.J. Boxma, R. Núñez-Queija, Heavy tails: the effect of the service discipline, in: T. Field, P.G. Harrison, J. Bradley, U. Harder (Eds.), *Computer Performance Evaluation-Modelling Techniques and Tools*, Proceedings of the Tools 2002, London, UK, Springer, Berlin, 2002, pp. 1–30.
- [17] S.C. Borst, O.J. Boxma, M.J.G. Van Uitert, Two coupled queues with heterogeneous traffic, in: J. Moreira de Souza, N.L.S. da Fonseca, E.A. de Souza e Silva (Eds.), *Teletraffic Engineering in the Internet Era*, Proceedings of the ITC-17, Salvador da Bahia, Brazil, North-Holland, Amsterdam, 2001, pp. 1003–1014.
- [18] S.C. Borst, O.J. Boxma, M.J.G. Van Uitert, The asymptotic workload behavior of two coupled queues, *Queue. Syst.* 43 (2003) 81–102.
- [19] S.C. Borst, A.P. Zwart, Fluid queues with heavy-tailed $M/G/\infty$ input, SPOR-Report 2001-02, Eindhoven University of Technology, submitted for publication.

- [20] O.J. Boxma, J.W. Cohen, The single server queue: heavy tails and heavy traffic, in: K. Park, W. Willinger (Eds.), *Self-Similar Network Traffic and Performance Evaluation*, Wiley, New York, 2000, pp. 143–169.
- [21] O.J. Boxma, J.W. Cohen, Q. Deng, Heavy-traffic analysis of the $M/G/1$ queue with priority classes, in: P. Key, D. Smith (Eds.), *Teletraffic Engineering in a Competitive World*, Proceedings of the ITC-16, Edinburgh, UK, North-Holland, Amsterdam, 1999, pp. 1157–1167.
- [22] O.J. Boxma, Q. Deng, J.A.C. Resing, Polling systems with regularly varying service and/or switchover times, *Adv. Perf. Anal.* 3 (2000) 71–107.
- [23] O.J. Boxma, V. Dumas, The busy-period in the fluid queue, *Perf. Eval. Rev.* 26 (1998) 100–110.
- [24] J. Cao, K. Ramanan, A Poisson limit for buffer overflow probabilities, in: *Proceedings of the IEEE Infocom*, New York, 2002, pp. 994–1003.
- [25] J.W. Cohen, Some results on regular variation for distributions in queueing and fluctuation theory, *J. Appl. Probab.* 10 (1973) 343–353.
- [26] J.W. Cohen, *The Single Server Queue*, revised edition, North-Holland, Amsterdam, 1982.
- [27] J.W. Cohen, O.J. Boxma, *Boundary Value Problems in Queueing System Analysis*, North-Holland, Amsterdam, 1983.
- [28] V.P. Chistyakov, A theorem on sums of independent, positive random variables and its applications to branching processes, *Theory Probab. Appl.* 9 (1964) 640–648.
- [29] D. Cline, Intermediate regular and Π variation, *Proc. London Math. Soc.* 68 (1994) 594–616.
- [30] M. Crovella, A. Bestavros, Self-similarity in World Wide Web traffic: evidence and possible causes, in: *Proceedings of the ACM Sigmetrics'96*, 1996, pp. 160–169.
- [31] Q. Deng, The two-queue $E/1 - L$ polling model with regularly varying service and/or switchover times, SPOR-Report 2001-09, Department of Mathematics and Computer Science, Eindhoven University of Technology; *Stoch. Mod.*, to appear.
- [32] G. Fayolle, R. Iasnogorodski, Two coupled processors: the reduction to a Riemann–Hilbert problem, *Z. Wahrsch. Verw. Gebiete* 47 (1979) 325–351.
- [33] W. Feller, *An Introduction to Probability Theory and its Applications*, vol. II, Wiley, New York, 1971.
- [34] P.R. Jelenković, P. Momčilović, Resource sharing with subexponential distributions, in: *Proceedings of the IEEE Infocom*, New York, 2002, pp. 1316–1325.
- [35] P.R. Jelenković, P. Momčilović, Large deviations of square-root insensitive random sums, Technical Report, Columbia University, 2002.
- [36] P.R. Jelenković, P. Momčilović, A.P. Zwart, Reduced-load equivalence under subexponentiality, Research Report RR 4444, INRIA Rocquencourt, 2002; *Queue. Syst.*, to appear.
- [37] J. Karamata, Sur un mode de croissance régulière des fonctions, *Mathematica Cluj* 4 (1930) 38–53.
- [38] L. Kleinrock, *Queueing Systems*, vol. II: Computer Applications, Wiley, New York, 1976.
- [39] C. Klüppelberg, Subexponential distributions and integrated tails, *J. Appl. Probab.* 25 (1988) 132–141.
- [40] A.G. Konheim, I. Meilijson, A. Melkman, Processor sharing of two parallel lines, *J. Appl. Probab.* 18 (1981) 952–956.
- [41] D.A. Korshunov, On distribution tail of the maximum of a random walk, *Stoch. Proc. Appl.* 72 (1997) 97–103.
- [42] W.E. Leland, M.S. Taqqu, W. Willinger, D.V. Wilson, On the self-similar nature of Ethernet traffic (extended version), *IEEE/ACM Trans. Netw.* 2 (1994) 1–15.
- [43] A. De Meyer, J.L. Teugels, On the asymptotic behaviour of the distributions of the busy period and the service time in $M/G/1$, *J. Appl. Probab.* 17 (1980) 802–813.
- [44] T. Mikosch, Regular variation, subexponentiality and their applications in probability theory, EURANDOM Report 99-013.
- [45] R. Núñez-Queija, Processor-sharing models for integrated-services networks, Ph.D. Thesis, Eindhoven University of Technology, ISBN 90-646-4667-8, 2000; also available from the author upon request.
- [46] R. Núñez-Queija, Sojourn-times in a processor-sharing queue with service interruptions, *Queue. Syst.* 34 (2000) 351–386.
- [47] R. Núñez-Queija, Queues with equally heavy sojourn-time and service requirement distributions, *Ann. Oper. Res.* 113 (2002) 101–117.
- [48] T.J. Ott, The sojourn-time distribution in the $M/G/1$ queue with processor sharing, *J. Appl. Probab.* 21 (1984) 360–378.
- [49] A.G. Pakes, On the tails of waiting-time distributions, *J. Appl. Probab.* 12 (1975) 555–564.

- [50] A.K. Parekh, R.G. Gallager, A generalized processor sharing approach to flow control in integrated services networks: the single-node case, *IEEE/ACM Trans. Netw.* 1 (1993) 344–357.
- [51] A. Paxson, S. Floyd, Wide area traffic: the failure of Poisson modeling, *IEEE/ACM Trans. Netw.* 3 (1995) 226–244.
- [52] S. Resnick, G. Samorodnitsky, Activity periods of an infinite server queue and performance of certain heavy-tailed fluid queues, *Queue. Syst.* 33 (1999) 43–71.
- [53] M. Sakata, S. Noguchi, J. Oizumi, An analysis of the $M/G/1$ queue under round-robin scheduling, *Oper. Res.* 19 (1971) 371–385.
- [54] L.E. Schrage, L.W. Miller, The queue $M/G/1$ with the shortest remaining processing time discipline, *Oper. Res.* 14 (1966) 670–684.
- [55] R. Schassberger, A new approach to the $M/G/1$ processor sharing queue, *Adv. Appl. Probab.* 16 (1984) 802–813.
- [56] W. Willinger, M.S. Taqqu, R. Sherman, D.V. Wilson, Self-similarity through high-variability: statistical analysis of Ethernet LAN traffic at the source level, *IEEE/ACM Trans. Netw.* 5 (1997) 71–86.
- [57] S.F. Yashkov, A derivation of response time distribution for a $M/G/1$ processor-sharing queue, *Probab. Control Inf. Theory* 12 (1983) 133–148.
- [58] S.F. Yashkov, Processor-sharing queues: Some progress in analysis, *Queue. Syst.* 2 (1987) 1–17.
- [59] A.P. Zwart, Sojourn-times in a multiclass processor sharing queue, in: P. Key, D. Smith (Eds.), *Teletraffic Engineering in a Competitive World*, Proceedings of the ITC-16, Edinburgh, UK, North-Holland, Amsterdam, 1999, pp. 335–344.
- [60] A.P. Zwart, Queueing systems with heavy tails, Ph.D. Thesis, Eindhoven University of Technology, 2001.
- [61] A.P. Zwart, Tail asymptotics for the busy period in the $GI/G/1$ queue, *Math. Oper. Res.* 26 (2001) 485–493.
- [62] A.P. Zwart, S.C. Borst, M. Mandjes, Exact asymptotics for fluid queues fed by multiple heavy-tailed On-Off flows, Shortened version in: Proceedings of the IEEE Infocom, 2001, Anchorage AK, USA, pp. 279–288, *Ann. Appl. Probab.*, in press.
- [63] A.P. Zwart, O.J. Boxma, Sojourn time asymptotics in the $M/G/1$ processor sharing queue, *Queue. Syst.* 35 (2000) 141–166.



S.C. Borst received the M.Sc. degree in applied mathematics from the University of Twente, The Netherlands, in 1990, and the Ph.D. degree from the University of Tilburg, The Netherlands, in 1994. During the fall of 1994, he was a visiting scholar at the Statistical Laboratory of the University of Cambridge, England. In 1995, he joined the Mathematics of Networks and Systems department of Bell Laboratories, Lucent Technologies in Murray Hill, USA. Since the fall of 1998, he has also been with the Center for Mathematics and Computer Science (CWI) in Amsterdam. He also has a part-time appointment as a professor of Stochastic Operations Research at Eindhoven University of Technology. He is a member of IFIP Working Group 7.3, and serves on the editorial board of several journals. His main research interests are in the performance evaluation of communication networks and computer systems.



O.J. Boxma (Ph.D. Utrecht, 1977) has been an IBM Postdoctoral Fellow during 1978–1979. He has worked at the University of Utrecht (1974–1985) and CWI (1985–1998). From 1987 until September 1998 he also was professor of Operations Research at Tilburg University. Since September 1998 he holds the chair of Stochastic Operations Research in Eindhoven University of Technology. In addition, he coordinates the Stochastic Networks program of the European Research Institute EURANDOM. Onno Boxma is co-author/co-editor of five books on queueing theory and performance evaluation. He serves on the editorial board of several journals, and he is a member of IFIP WG7.3 and of the International Advisory Board of the ITC. His main research interests are in queueing theory and its applications to the performance analysis of computer-communication and production systems.



R. Núñez-Queija received his master's degree in econometrics from the Econometrics Department of the Free University of Amsterdam in 1995 and a Ph.D. from the Mathematics and Computer Science Department of Eindhoven University of Technology in 2000. He was a post-doc at INRIA (Sophia Antipolis, France) in 2000. Currently he is a member of the Probability, Networks and Algorithms department at CWI (Center for Mathematics and Computer Science, Amsterdam) and assistant professor in Stochastic Operations Research at the Faculty of Mathematics and Computer Science of Eindhoven University of Technology. His main research interests are in queueing theory and the performance analysis of communication networks.



A.P. Zwart got his Ph.D. at Eindhoven University of Technology on September 11, 2001. After that, he spent one year as postdoctoral fellow at INRIA (Rocquencourt, France). Currently, he is assistant professor at the Department of Mathematics and Computer Science, again in Eindhoven. His research is funded by the Dutch Research program for innovational research.