



**REPORT** *RAPPORT*

*PNA*

Probability, Networks and Algorithms



*Probability, Networks and Algorithms*

Sojourn times in the  $M/G/1$  FB queue with light-tailed service times

M.R.H. Mandjes, M. Nuyens

**REPORT PNA-E0411 JUNE 2004**

CWI is the National Research Institute for Mathematics and Computer Science. It is sponsored by the Netherlands Organization for Scientific Research (NWO).

CWI is a founding member of ERCIM, the European Research Consortium for Informatics and Mathematics.

CWI's research has a theme-oriented structure and is grouped into four clusters. Listed below are the names of the clusters and in parentheses their acronyms.

**Probability, Networks and Algorithms (PNA)**

Software Engineering (SEN)

Modelling, Analysis and Simulation (MAS)

Information Systems (INS)

Copyright © 2004, Stichting Centrum voor Wiskunde en Informatica

P.O. Box 94079, 1090 GB Amsterdam (NL)

Kruislaan 413, 1098 SJ Amsterdam (NL)

Telephone +31 20 592 9333

Telefax +31 20 592 4199

ISSN 1386-3711

# Sojourn times in the M/G/1 FB queue with light-tailed service times

## ABSTRACT

The asymptotic decay rate of the sojourn time of a customer in the stationary M/G/1 queue under the Foreground-Background (FB) service discipline is studied. The FB discipline gives service to those customers that have received the least service so far. We prove that for light-tailed service times the decay rate of the sojourn time is equal to the decay rate of the busy period. It is shown that FB minimises the decay rate in the class of work-conserving disciplines.

*2000 Mathematics Subject Classification:* Primary 60K25, Secondary 68M20; 90B22

*Keywords and Phrases:* decay rate, sojourn time, Foreground-Background (FB), LAST, service discipline, light tails, busy period

# Sojourn times in the M/G/1 FB queue with light-tailed service times

M. Mandjes\*      M. Nuyens†

**Keywords:** decay rate, sojourn time, Foreground-Background (FB), LAST, service discipline, light tails, busy period

**AMS 2000 Subject Classification:** Primary 60K25, Secondary 68M20; 90B22

## Abstract

The asymptotic decay rate of the sojourn time of a customer in the stationary M/G/1 queue under the Foreground-Background (FB) service discipline is studied. The FB discipline gives service to those customers that have received the least service so far. We prove that for light-tailed service times the decay rate of the sojourn time is equal to the decay rate of the busy period. It is shown that FB minimises the decay rate in the class of work-conserving disciplines.

## 1 Introduction

The sojourn time of a customer, i.e. the time between his arrival and departure, is an often used performance measure for queues. In this note we compute the asymptotic decay rate of the tail of the sojourn-time distribution of the stationary M/G/1 queue with the Foreground-Background (FB) discipline. This decay rate is then used to compare the performance of FB with other service disciplines like PS and FIFO.

---

\*CWI, Amsterdam, The Netherlands, and University of Twente, Faculty of Mathematical Sciences, The Netherlands

†KdV Institute for Mathematics, University of Amsterdam, The Netherlands

The FB discipline gives service to those customers who have received the least amount of service so far. If there are  $n$  such customers, each of them is served at rate  $1/n$ . Thus, when the *age* of a customer is the amount of service a customer has received, the FB discipline gives priority to the *youngest* customers. In the literature this discipline has been called LAS or LAST (least-attained service time first) as well.

Let  $V$  denote the sojourn time of a customer in the stationary M/G/1 FB queue. Núñez Queija [6] showed that for service-time distributions with regularly varying tails of index  $\eta \in (1, 2)$ , the distribution of  $V$  satisfies

$$P(V > x) \sim P(B > (1 - \rho)x), \quad x \rightarrow \infty, \quad (1)$$

where  $\rho$  is the load of the system,  $B$  is the generic service time, and  $\sim$  means that the quotient converges to 1. Using Núñez Queija's method, Nuyens [7] obtained (1) under weaker assumptions. In case of regularly varying service times the tail of  $V$  under other disciplines, like FIFO, LIFO, PS and SRPT, has been found to be heavier than under FB, see Borst, Boxma, Núñez Queija and Zwart [1].

Additional support for the effective performance of FB under heavy tails is given by Righter and Shanthikumar [8, 9, 10]. They show that for certain classes of service times (including e.g. the Pareto distribution), the FB discipline minimises the queue length, measured in number of customers, in the class of all disciplines that do not know the exact value of the service times.

For light-tailed service times the FB discipline does not perform so well, although for gamma densities  $\lambda^\alpha x^{\alpha-1} \exp(-\lambda x) / \Gamma(\alpha)$  with  $0 < \alpha \leq 1$ , FB still minimises the queue length, and for exponential service times the queue length is independent of the service discipline. However, for many other light-tailed service times, for example those with a decreasing failure rate, the queue shows opposite behaviour and the queue length is maximised by FB, see Righter and Shanthikumar [8, 9, 10]. This undesirable behaviour of the FB discipline is very pronounced for deterministic service times. In this extreme case in the FB queue all customers stay till the end of the busy period, and the sojourn time under the FB discipline is *maximal* in the class of all work-conserving disciplines. In this note we consider the (asymptotic) decay rate of the sojourn time, where the (*asymptotic*) *decay rate*  $dr(X)$  of a random variable  $X$  is defined as

$$dr(X) = \left| \lim_{x \rightarrow \infty} x^{-1} \log P(X > x) \right|,$$

given that the limit exists. Hence a larger decay rate means a smaller probability that the random variable takes on very large values. In this sense sojourn times are better when they have larger decay rates.

It turns out that for the M/G/1 FB queue in which the service-time distribution has an exponentially fast decreasing tail, large sojourn times are relatively likely, in the following sense. Assume that the service times have a finite exponential moment, or equivalently, the Laplace transform is analytic in a neighbourhood of zero. The main theorem of this note is then the following.

**Theorem 1** *Let  $V$  be the sojourn time of a customer in the stationary M/G/1 FB queue, and let  $L$  be the length of a busy period. If the service-time distribution has a finite exponential moment, then the decay rate of  $V$  exists and satisfies*

$$dr(V) = dr(L). \tag{2}$$

It is shown below that the decay rate of the sojourn time in an M/G/1 queue with any work-conserving discipline is bounded from below by the decay rate of the residual life of a busy period. For service times with an exponential moment the latter decay rate is equal to that of a normal busy period. Hence (2) is the lowest possible decay rate for the sojourn time under a work-conserving discipline. Using the decay rate of  $V$  as a criterion to measure the performance of a service discipline then leads to the following conclusion: for service times with an exponential moment, the FB discipline is the worst discipline in the class of work-conserving disciplines.

The paper is organised as follows. In Section 2 we present the notation, some preliminaries, and prove the lower bound for the decay rate of the sojourn time under any work-conserving discipline. In Section 3 Theorem 1 is proved. Section 4 discusses the result and the decay rate of the sojourn time in queues operating under several other service disciplines.

## 2 Preliminaries

Throughout this note we assume that the generic service time  $B$  with distribution function  $F$  in the M/G/1 queue satisfies the following assumption.

**Assumption 1** *The generic service time  $B$  has an exponential moment, i.e.,*

$$E \exp(\gamma B) < \infty$$

*for some  $\gamma > 0$ .*

Let in addition the stability condition  $\rho = \lambda EB < 1$  hold, where  $\lambda$  is the rate of the Poisson arrival process. The proofs in this note rely on some properties of the busy-period length  $L$  and related random variables, which we derive in this section.

Under assumption 1, Cox and Smith [3] have shown that  $P(L > x) \sim bx^{-3/2}e^{-cx}$  for certain constants  $b, c > 0$ . In particular,  $L$  has decay rate  $c$ . In fact, by expression (46) on page 154 of Cox and Smith [3],  $c = \lambda - \zeta - \lambda g(\zeta)$ , where  $g$  is the Laplace transform of the service-time distribution, and  $\zeta < 0$  is such that  $g'(\zeta) = -\lambda^{-1}$ . Hence  $\zeta$  is the root of the derivative of the function  $m(x) = \lambda - x - \lambda g(x)$ . Since  $m(x)$  attains its maximum in the point  $\zeta$ , we may write  $c$  in terms of the Legendre transform of  $B$ ,

$$c = dr(L) = \sup_{\theta} \{\theta - \lambda(Ee^{\theta B} - 1)\}. \quad (3)$$

**Remark** This expression shows up as well in the following context. Consider a Poisson stream, with intensity  $\lambda$ , of i.i.d. jobs, where every job is distributed according to the random variable  $B$ . Let  $A(x)$  denote the amount of work generated in an arbitrary time window of length  $x$ . It is an easy corollary of Cramér's theorem that

$$\lim_{x \rightarrow \infty} \frac{1}{x} \log P(A(x) > x) = -\sup_{\theta} \{\theta - \log Ee^{\theta A(1)}\}. \quad (4)$$

Noting that

$$Ee^{\theta A(1)} = \sum_{k=0}^{\infty} e^{-\lambda} \frac{\lambda^k}{k!} (Ee^{\theta B})^k = \exp(\lambda(Ee^{\theta B} - 1)),$$

we observe that  $P(L > x)$  and  $P(A(x) > x)$  have the same decay rate. This is somewhat surprising, as  $\{A(x) > x\}$  obviously depends just on  $A(x)$ , i.e. the amount of traffic in a window of length  $x$ , whereas  $\{L > x\}$  depends on  $A(y)$  for *all*  $y \in [0, x]$ , due to

$$\{L > x\} \stackrel{d}{=} \{B_1 + A(y) > y, \forall y \in [0, x]\}.$$

Here  $B_1$  is the first service time in the busy period  $L$ .

In renewal theory the notion of *residual life*, also known as excess or forward-recurrence time, is standard. Let  $\tilde{L}$  be the residual life of a busy period. Then  $P(\tilde{L} > x) = (EL)^{-1} \int_x^{\infty} P(L > y) dy$ , see for instance Cox [2]. Using standard calculus we find

$$dr(\tilde{L}) = \left| \lim_{x \rightarrow \infty} \frac{1}{x} \log \int_x^{\infty} y^{-3/2} e^{-cy} dy \right| = c = dr(L). \quad (5)$$

Hence  $\tilde{L}$  has the same decay rate as  $L$ .

Another ingredient used in the proofs below is the M/G/1 queue with truncated

generic service time  $B \wedge \tau$ ,  $\tau > 0$ . Call this the  $\tau$ -queue and let  $L(\tau)$  denote the length of a busy period (a  $\tau$ -busy period) in this queue. Let  $\tilde{L}(\tau)$  be its the residual life and define  $L^*(\tau)$  to be the length of a  $\tau$ -busy period in which the first service time  $B_1$  is at least  $\tau$ , i.e.

$$P(L^*(\tau) > x) = P(L(\tau) > x \mid B_1 \geq \tau).$$

We now show that the random variables  $L(\tau)$ ,  $\tilde{L}(\tau)$  and  $L^*(\tau)$  have the same decay rate.

**Lemma 2** *Let  $\tau > 0$  be such that  $P(B \geq \tau) > 0$ . Then*

$$dr(L(\tau)) = dr(L^*(\tau)) = dr(\tilde{L}(\tau)) > 0.$$

**Proof** We show that  $L$  and  $L^*$  have the same decay rate. The proof is then finished by using (5). Let  $B_1$  denote the first service time in the busy period, hence  $B_1 \stackrel{d}{=} B$ . Assume that  $\tau > 0$  is such that  $P(B \geq \tau) > 0$ . If  $B_1 \geq \tau$ , then the first service time is maximal in the  $\tau$ -queue, as all service times are bounded by  $\tau$ . Hence

$$P(L(\tau) > x) \leq P(L(\tau) > x \mid B_1 \geq \tau) = P(L^*(\tau) > x), \quad x \geq 0.$$

Further,

$$\begin{aligned} P(L(\tau) > x) &\geq P(L(\tau) > x, B_1 \geq \tau) = P(L(\tau) > x \mid B_1 \geq \tau)P(B_1 \geq \tau) \\ &= P(L^*(\tau) > x)P(B_1 \geq \tau), \end{aligned} \tag{6}$$

From (6) it follows that  $P(L^*(\tau) > x)$  and  $P(L(\tau) > x)$  differ only by a term independent of  $x$ . Hence  $dr(L) = dr(\tilde{L})$ . ■

In this note we need the following lemma about the decay rate of the sum of two independent random variables.

**Lemma 3** *Let  $X$  and  $Y$  be non-negative, independent random variables such that  $dr(X) = dr(Y) = \alpha$  for some  $\alpha > 0$ . Then also  $dr(X + Y) = \alpha$ .*

**Proof** Since both  $X$  and  $Y$  are positive,  $-\alpha$  is clearly a lower bound for  $dr(X + Y)$ . For the upper bound let  $n \in \mathbb{N}$  be fixed. Then,

$$P(X + Y > x) \leq \sum_{i=0}^{n-1} P\left(X \geq \frac{ix}{n}\right)P\left(Y \geq \frac{(n-i-1)x}{n}\right).$$



Fix  $\varepsilon > 0$ . For  $x$  sufficiently large, for all  $i \in \{0, \dots, n-1\}$ ,

$$\begin{aligned} P\left(X \geq \frac{ix}{n}\right)P\left(Y \geq \frac{(n-i-1)x}{n}\right) &\leq \exp\left(-(\alpha-\varepsilon)\frac{ix}{n} - (\alpha-\varepsilon)x\frac{n-i-1}{n}\right) \\ &= \exp\left(-(\alpha-\varepsilon)\frac{(n-1)x}{n}\right). \end{aligned}$$

Hence,

$$\limsup_{x \rightarrow \infty} \frac{1}{x} \log P(X+Y > x) \leq -(\alpha-\varepsilon)\left(1 - \frac{1}{n}\right). \quad (7)$$

Since (7) holds for every  $n \in \mathbb{N}$  and  $\varepsilon > 0$ , we may take the limits  $n \rightarrow \infty$  and  $\varepsilon \downarrow 0$ , and the result follows.  $\blacksquare$

Let  $D$  be the time from the arrival of a customer till the first moment that the system is empty. The following proposition is valid also in the case that Assumption 1 does not hold.

**Proposition 4** *Consider a stationary queue with an arbitrary service-time distribution, Poisson arrivals and a work-conserving discipline. Then  $D \stackrel{d}{=} A\tilde{L} + L$ , where  $P(A=1) = \rho = 1 - P(A=0)$  and  $A, \tilde{L}$  and  $L$  are independent.*

**Proof** The value of the random variable  $D$  does not depend on the service discipline. There are two possibilities. With probability  $1 - \rho$  the customer finds the system empty. In this case  $D$  is just the length  $L$  of the busy period started by the customer. Secondly, if our customer enters a busy system, then the server may first finish all the work in the system apart from the work of our tagged customer. The moment the remainder of the original busy period, which has length  $\tilde{L}$ , is finished, our customer starts a sub-busy period. This length of this sub-busy period, which is independent of  $\tilde{L}$ , is distributed like  $L$ .  $\blacksquare$

For the stationary  $\tau$ -queue with Poisson arrivals and a work-conserving discipline, we have the following corollary.

**Corollary 5** *In the stationary  $\tau$ -queue, the random variable  $D$  satisfies*

$$D \stackrel{d}{=} A(\tau)\tilde{L}(\tau) + L(\tau),$$

where  $P(A(\tau) = 1) = \lambda E(B \wedge \tau)$ . If the customer has service time  $\tau$  in the  $\tau$ -queue, then  $D \stackrel{d}{=} A(\tau)\tilde{L}(\tau) + L^*(\tau)$ .

Since the system is work-conserving, the sojourn time of a customer is not longer than  $D$ . Hence  $V \leq_{st} D$  for every service discipline. Since  $A\tilde{L}$  and  $L$  satisfy the conditions of Lemma 3, the following corollary holds.

**Corollary 6** *For every work-conserving service discipline, the sojourn time  $V$  of a customer in the stationary queue satisfies*

$$\limsup_{x \rightarrow \infty} \frac{1}{x} \log P(V > x) \leq \lim_{x \rightarrow \infty} \frac{1}{x} \log P(A\tilde{L} + L > x) = -dr(L).$$

An immediate consequence of this Corollary and Theorem 1, which will be proved in the next section, is the following.

**Corollary 7** *The FB discipline minimises the decay rate of the sojourn time in the class of work-conserving disciplines.*

In Section 4 it is discussed that there are service disciplines with a strictly larger decay rate, e.g. FIFO.

Interestingly, for service times with certain Gamma distributions, the FB discipline minimises the queue length, as was mentioned in the introduction, but the sojourn time has the smallest decay rate. This shows that optimising one characteristic in a queue may have an ill effect on other characteristics.

The existence of a finite exponential moment in the corollary is crucial: for heavy-tailed service times the tail of  $V$  cannot be bounded by that of  $L$ . For example, in the M/G/1 FIFO queue with service times satisfying  $P(B > x) = x^{-\nu} \mathcal{L}(x)$ , where  $\mathcal{L}(x)$  is a slowly varying function at  $\infty$  and  $\nu > 1$ , De Meyer and Teugels [4] showed that

$$P(L > x) \sim (1 - \rho)^{-\nu-1} x^{-\nu} \mathcal{L}(x).$$

It may be seen that in this case the tail of  $\tilde{B}$ , the residual life of the generic service time  $B$ , is one degree heavier than that of  $B$ . Now note that for the FIFO discipline we have  $V_{\text{FIFO}} \geq A\tilde{B}$ . Hence the tail of  $V$  is at least one degree heavier than that of  $L$ , see also Borst *et al.* [1] for further references. In the light-tailed case this phenomenon is absent since the tails of  $L$  and  $\tilde{L}$  have the same decay rate.

### 3 Proof of the theorem

In this section Theorem 1 is proved. The results in this section rely on the following decomposition of  $V$ . Let  $V(\tau)$  be the sojourn time in the stationary M/G/1 queue of

a customer with service time  $\tau$ . The sojourn time  $V$  of an arbitrary customer in the stationary queue satisfies

$$P(V > x) = \int P(V(\tau) > x) dF(\tau). \quad (8)$$

Here  $F$  is the service-time distribution. Hence we may write

$$P(V > x) = E_B P(V(B) > x),$$

where  $B$  is a generic service time independent of  $V(\tau)$ , and  $E_B$  denotes the expectation w.r.t.  $B$ . Theorem 1 is proved using this representation of  $V$ . In the next lemma we compute the decay rate of  $V(\tau)$ .

**Proposition 8** *Let  $\tau > 0$  be such that  $P(B \geq \tau) > 0$ . If the service-time distribution satisfies Assumption 1, then  $dr(V(\tau)) = dr(L(\tau))$ .*

**Proof** By the nature of the FB discipline, the sojourn time  $V(\tau)$  of a customer with service time  $\tau$  who enters a stationary queue is the time till the first epoch that no customers younger than  $\tau$  are present. This is the time till the end of the  $\tau$ -busy period that he either finds in the  $\tau$ -queue, or starts. By Corollary 5,  $V(\tau)$  then satisfies

$$V(\tau) \stackrel{d}{=} A(\tau)\tilde{L}(\tau) + L^*(\tau), \quad (9)$$

where  $\tilde{L}(\tau)$  is the residual life of a  $\tau$ -busy period,  $L^*(\tau)$  is a  $\tau$ -busy period that starts with a customer with service time  $\tau$ ,  $P(A(\tau) = 1) = 1 - P(A(\tau) = 0) = \lambda E(B \wedge \tau)$  and  $A(\tau)$ ,  $\tilde{L}(\tau)$  and  $L^*(\tau)$  are independent. By Lemma 2 the random variables  $A(\tau)\tilde{L}(\tau)$  and  $L^*(\tau)$  satisfy the condition of Lemma 3. From (9) and again Lemma 2, it follows that

$$dr(V(\tau)) = dr(A(\tau)\tilde{L}(\tau) + L^*(\tau)) = dr(L(\tau)). \quad (10)$$

This completes the proof. ■

Having found the upper bound for the decay rate in Corollary 6, the following lemma provides the basis for finding the lower bound. The *endpoint*  $x_F$  of the service-time distribution  $F$  is defined as  $x_F = \inf\{u \geq 0 : F(u) = 1\}$ .

**Lemma 9** *Let  $V$  be the sojourn time of a customer in the stationary M/G/1 FB queue. Suppose the service-time distribution satisfies Assumption 1. If  $\tau_0 > 0$  and  $P(B \geq \tau_0) > 0$ , then*

$$\liminf_{x \rightarrow \infty} \frac{1}{x} \log P(V > x) \geq -P(B \geq \tau_0)^{-1} \int_{[\tau_0, x_F]} dr(L(\tau)) dF(\tau). \quad (11)$$

Here  $F$  is the distribution function of the generic service time  $B$ .

**Proof** Let  $B$  and  $V$  denote the service time and the sojourn time of a customer in the stationary queue. Let  $\tau_0 > 0$  be such that  $P(B \geq \tau_0) > 0$ . Then

$$P(V > x) \geq P(V > x, B \geq \tau_0) = P(V > x \mid B \geq \tau_0)P(B \geq \tau_0). \quad (12)$$

Using the representation (8), we find

$$\log P(V > x \mid B \geq \tau_0) = \log E_B[P(V(B) > x) \mid B \geq \tau_0]. \quad (13)$$

Since  $\log x$  is a concave function, applying Jensen's inequality to the conditional expectation in (13) yields

$$\log E_B[P(V(B) > x) \mid B \geq \tau_0] \geq E_B[\log P(V(B) > x) \mid B \geq \tau_0]. \quad (14)$$

From (12), (13) and (14) it follows that  $\Theta := \liminf_{x \rightarrow \infty} \frac{1}{x} \log P(V > x)$  satisfies

$$\Theta \geq \liminf_{x \rightarrow \infty} \frac{1}{x} \log \int_{[\tau_0, x_F]} \log P(V(\tau) > x) dF(\tau) / P(B \geq \tau_0). \quad (15)$$

Applying Fatou's lemma to (15) yields

$$\Theta \geq P(B \geq \tau_0)^{-1} \int_{[\tau_0, x_F]} \lim_{x \rightarrow \infty} \frac{1}{x} \log P(V(\tau) > x) dF(\tau).$$

The result now follows from Proposition 8. ■

The following lemma is used to develop the lower bound for the decay rate of  $V$  from Lemma 9. We introduce the notation  $c(\tau) = dr(L(\tau))$ , so that  $c = dr(L) = c(x_F)$ .

**Lemma 10** *The function  $c(\tau)$  is decreasing in  $\tau$ . Furthermore,  $c(\tau) \rightarrow c(x_F)$  as  $\tau \rightarrow x_F$ .*

**Proof** For all  $\tau$ , the function  $h_\tau(\theta) = \theta - \lambda(Ee^{\theta(B \wedge \tau)} - 1)$  is concave in  $\theta$ , since any moment generating function is convex. Furthermore  $\lim_{\theta \rightarrow -\infty} h_\tau(\theta) = \lim_{\theta \rightarrow \infty} h_\tau(\theta) = -\infty$ . By definition of  $L(\tau)$  and (3), we may write  $c(\tau) = \sup_\theta \{h_\tau(\theta)\}$ . Then  $c(\tau)$  is decreasing in  $\tau$ , since  $h_\tau(\theta)$  is decreasing in  $\tau$ . Since  $c(\tau) \geq h_\tau(0) = 0$  for all  $\tau$ , and  $c(\tau)$  is decreasing,  $c(\tau)$  converges for  $\tau \rightarrow x_F$ . Now note that  $h_\tau(\theta)$  is continuous in  $\tau$  for all  $\theta \in [0, \sup\{\eta : Ee^{\eta B} < \infty\})$ , even if  $B$  has a discrete distribution. Since the supremum of  $\theta - \lambda(Ee^{\theta B} - 1)$  is attained in this interval, we have  $\lim_{\tau \rightarrow x_F} c(\tau) = c(x_F)$ . ■

**Proposition 11** *Let  $V$  be the sojourn time of a customer in the stationary  $M/G/1$  FB queue. If the service-time distribution satisfies Assumption 1, then*

$$\liminf_{x \rightarrow \infty} \frac{1}{x} \log P(V > x) \geq -dr(L). \quad (16)$$

**Proof** If  $P(B = x_F) > 0$ , then choosing  $\tau_0 = x_F$  in (11) yields

$$\liminf_{x \rightarrow \infty} \frac{1}{x} \log P(V > x) \geq -c(x_F) = -dr(L),$$

and (16) holds. Assume  $P(B = x_F) = 0$ , and let  $\varepsilon > 0$ . By Lemma 10 there exists an  $x_\varepsilon < x_F$  such that  $c(\tau) \leq c + \varepsilon$  for all  $\tau \geq x_\varepsilon$ . Choosing  $\tau_0 = x_\varepsilon$  in (11) then yields

$$\begin{aligned} \liminf_{x \rightarrow \infty} \frac{1}{x} \log P(V > x) &\geq -P(B \geq x_\varepsilon)^{-1} \int_{[x_\varepsilon, x_F]} c(\tau) dF(\tau) \\ &\geq -P(B \geq x_\varepsilon)^{-1} \int_{[x_\varepsilon, x_F]} (c + \varepsilon) dF(\tau) = -c - \varepsilon. \end{aligned}$$

Since  $\varepsilon > 0$  was arbitrary, the lower bound (16) follows. ■

**Proof of Theorem 1** The upper bound is established in Corollary 6 and the lower bound in Proposition 11. ■

## 4 Discussion

The decay rate of the sojourn time  $V$  in the M/G/1 FB queue is the same as for the preemptive LIFO queue. Indeed, the sojourn time of a customer in the stationary M/G/1 queue under the preemptive LIFO discipline is just the length of the sub-busy period started by that customer. From Theorem 1 it follows that the decay rates of the sojourn times for LIFO and FB are equal.

The sojourn time of a customer in the stationary queue under FIFO satisfies  $V_{\text{FIFO}} = B + W$ , where  $W$  is the stationary workload. From the Pollaczek-Khinchin formula,

$$Ee^{-sW} = \frac{s(1 - \rho)}{s - \lambda + \lambda E \exp(-sB)}, \quad (17)$$

it follows that the decay rate of  $W$  is the value of  $s$  for which the denominator in (17) vanishes. Hence  $dr(W)$  is the positive root  $\theta_0$  of  $h(\theta) = \theta - \lambda(Ee^{\theta B} - 1)$ . Furthermore, since  $dr(B) = \inf\{\theta : h(\theta) = -\infty\}$ , we have  $\theta_0 < dr(B) \leq \infty$ . An analogue of Lemma 3 then yields that  $c_{\text{FIFO}} := dr(V_{\text{FIFO}}) = \theta_0$ .

Since  $h$  is concave,  $h(0) = 0$  and  $h'(0) = 1 - \lambda EB < 1$ , we have by Theorem 1 and (3) that

$$c_{\text{FB}} := dr(V_{\text{FB}}) = dr(L) = \sup_{\theta} h(\theta) < \theta_0 = c_{\text{FIFO}} < dr(B), \quad (18)$$

see also Figure 2 below. Hence, in the FIFO system, the decay rate of the sojourn time is strictly larger than that in the FB queue. As an illustration, consider the M/M/1

queue in which the service times have expectation  $1/\mu$ . For stability we assume  $\lambda < \mu$ . Straightforward computations then yield that  $c_{\text{FB}} = (\sqrt{\mu} - \sqrt{\lambda})^2$ ,  $c_{\text{FIFO}} = \mu - \lambda$  and  $dr(B) = \mu$ . Since  $\lambda < \mu$ , we conclude that for the M/M/1 queue, inequality (18) is satisfied.

Finally, Mandjes en Zwart [5] consider the PS queue with light-tailed service requests. They show that the decay rate of  $P(V_{\text{PS}} > x)$  is equal to  $dr(L)$  as well, under the additional requirement that, for any positive constant  $k$ ,

$$\lim_{x \rightarrow \infty} \frac{1}{x} \log P(B > k \log x) = 0.$$

For deterministic requests, clearly this criterion is not met. Indeed, in [5] it is shown that the decay rate of  $V$  in the M/D/1 queue with the PS discipline is larger than  $dr(L)$ .

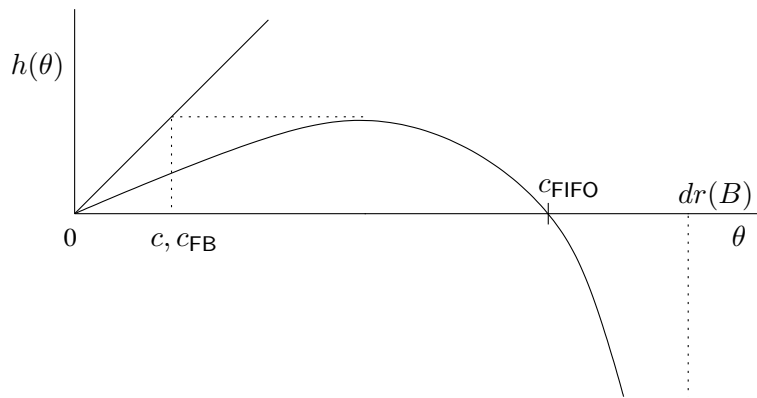


Figure 1 *The decay rates of the sojourn time under FB and FIFO.*

**Acknowledgements** The authors thank A.P. Zwart for kindly commenting on an earlier version of the paper, and the referee for his clear suggestions and comments, which have seriously improved the presentation of the paper.

## References

- [1] BORST, S., BOXMA, O., AND NÚÑEZ QUELJA, R. Heavy tails: the effect of the service discipline. In *Computer Performance Evaluation - Modelling Techniques and Tools. Proc. TOOLS 2002* (2002), T. Field, P. Harrison, J. Bradley, and U. Harder, Eds., pp. 1–30.

- [2] COX, D. *Renewal theory*. Methuen, 1962.
- [3] COX, D., AND SMITH, W. *Queues*. Methuen, 1961.
- [4] DE MEYER, A., AND TEUGELS, J. On the asymptotic behaviour of the distributions of the busy period and service time in M/G/1. *Journal of Applied Probability* 17, 3 (1980), 802–813.
- [5] MANDJES, M., AND ZWART, B. Large deviations for waiting times in processor sharing queues. *Submitted*.
- [6] NÚÑEZ QUEIJA, R. *Processor-Sharing Models for Integrated-Services Networks*. PhD thesis, Eindhoven University, 2000.
- [7] NUYENS, M. *Analysis of the Foreground Background queue*. PhD thesis, University of Amsterdam, to appear in 2004.
- [8] RIGHTER, R. Scheduling. In *Stochastic orders and their applications*, M. Shaked and R. Shanthikumar, Eds. Academic Press, 1994.
- [9] RIGHTER, R., AND SHANTHIKUMAR, J. Scheduling multiclass single server queueing systems to stochastically maximize the number of successful departures. *Probability in the Engineering and Informational Sciences* 3 (1989), 323–333.
- [10] RIGHTER, R., SHANTHIKUMAR, J., AND YAMAZAKI, G. On extremal service disciplines in single-stage queueing systems. *Journal of Applied Probability* 27 (1990), 409–416.