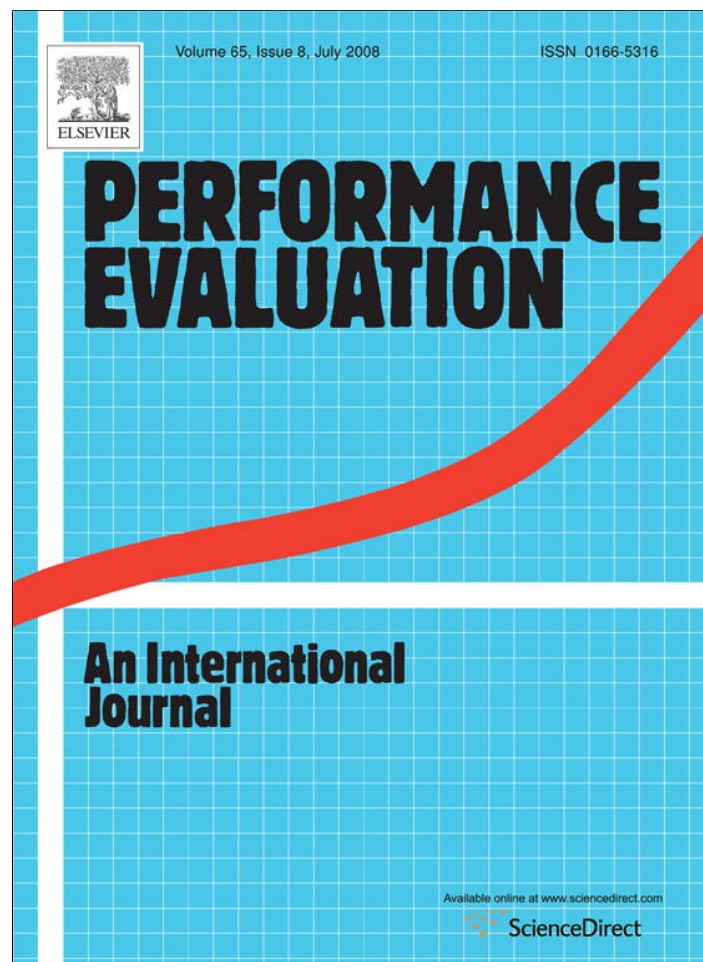


Provided for non-commercial research and education use.
Not for reproduction, distribution or commercial use.



This article appeared in a journal published by Elsevier. The attached copy is furnished to the author for internal non-commercial research and education use, including for instruction at the authors institution and sharing with colleagues.

Other uses, including reproduction and distribution, or selling or licensing copies, or posting to personal, institutional or third party websites are prohibited.

In most cases authors are permitted to post their version of the article (e.g. in Word or Tex form) to their personal website or institutional repository. Authors requiring further information regarding Elsevier's archiving and manuscript policies are encouraged to visit:

<http://www.elsevier.com/copyright>



Buffer overflow asymptotics for multiplexed regulated traffic

Y. Ying^a, F. Guillemin^b, R. Mazumdar^{c,*}, C. Rosenberg^c

^a School of ECE, Purdue University, West Lafayette, IN 47906, USA

^b France Telecom R&D F-22300, Lannion, France

^c School of ECE, University of Waterloo, Waterloo, Ontario, Canada, N2L 3G1

Received 13 October 2005; received in revised form 15 October 2007; accepted 16 October 2007

Available online 28 October 2007

Abstract

By characterizing the worst-case profile, which maximizes the content of a buffer fed with leaky bucket regulated flows in packet telecommunication networks, we derive a tight upper bound in the many-sources regime for the tail distribution of the workload generated by these flows in a FIFO queue with constant service rate. Furthermore, we compare this workload distribution with an $M/G/1$ queue and get insights on the better-than-Poisson property of regulated flows. We conclude that the superposition of independent regulated flows generates an asymptotically smaller workload than a marked Poisson process whose service times and intensity depend on the parameters of regulated sources.

© 2007 Elsevier B.V. All rights reserved.

Keywords: Network calculus; Statistical multiplexing; Better than Poisson

1. Introduction

In order to gain efficiency in bandwidth utilization in packet networks, standardization bodies like the IETF propose the characterization flows of information not only by a single parameter describing the peak information rate but also by means of additional parameters reflecting the random fluctuations of a flow. One common method for characterizing flows is to regulate their output profiles to conform to a specified envelope. The leaky bucket algorithm is one such method whereby, in addition to the peak rate π , the long-term mean rate ρ and the bucket size σ that is a measure of its burstiness are specified. A source conforming to the parameters (σ, ρ, π) is said to be (σ, ρ, π) -regulated.

A formalism to study the performance of a network supporting regulated flows, called *network calculus*, has been developed by Cruz [11,12] and more recently by Le Boudec [5] and Chang [7]. Network calculus provides a framework to compute end-to-end worst-case delay based on properties of the $(\min, +)$ algebra. However, it is essentially a deterministic approach which gives conservative worst-case bounds for network resource allocation.

In this paper, we exploit the statistical multiplexing features of regulated flows to obtain tighter performance bounds at a single node, keeping in mind that the ultimate goal is to derive end-to-end bounds. In particular, we address the

* Corresponding address: Department of Electrical and Computer Engineering, School of ECE, University of Waterloo, Waterloo, Ontario, Canada, N2L 3G1. Tel.: +1 519 888 4567; fax: +1 519 746 3077.

E-mail address: mazum@ece.uwaterloo.ca (R. Mazumdar).

issue of estimating the workload distribution of a FIFO queue fed with independent regulated sources. Our approach is based on asymptotic analysis of the workload distribution.

In order to provide satisfactory statistical guarantees, the workload distribution of a queue with (σ, ρ, π) -regulated flows has been extensively studied. These have focused on the overflow probability and/or delay distribution for an arbitrary number of input flows. Kesidis and Konstantopoulos [22,23] studied the problem via characterizing the extremal traffic shape, which maximizes the fraction of time when the buffer content is above a threshold b . Chang et al. [8] also derived an upper bound of the tail distribution of the workload but via partitioning the busy period.

As pointed out by the authors of [32] (see also [15]), these approaches can be regarded as applications of the Hoeffding bound given in [21], and they obtain the upper bounds via majorizing the queue length by a sum of independent processes and computing bounds based on these independent processes. Another approach to studying the superposition of regulated sources in a buffer has been presented by Busson and Massoulié [29], where they used Hoeffding's inequality based on the fact that the total number of packet arrivals in a given time interval from a regulated source is bounded. This approach was recently extended by Vojnovic and Le Boudec [32]. They also consider the so-called many-sources asymptotic framework, which is a key regime when the peak rate of sources is small compared with the server capacity and there are many sources; other results on many-sources asymptotics can be found in [4,8,10,26]. See also [31] for generalizations and alternative proofs of the results in [8,22,23].

In this paper we consider the problem of estimating buffer overflows in two different contexts of interest. The first is the many-sources asymptotic mentioned above, but our approach is different in that we first characterize the extremal traffic profile that maximizes the probability of buffer level exceedance. The second issue we consider is that of deriving simple majorization results for such flows that are valid even when the many-sources context is no more valid. This is with a view of characterizing the general properties of regulated flows that might be useful to obtain buffer overflow estimates inside the network. This, however, turns out to be a challenging problem since scheduling and multiplexing alter the initial statistical properties of traffic flows and increase the dependence between flows sharing common queues. This problem has been addressed by many authors; see for instance [3,14,9,28,33,34]. In particular, Bonald et al. in [3] compare the workload distribution of general packet arrival processes with that of an $M/G/1$ queue with Poisson/MTU arrivals (i.e., when each packet has fixed size equal to MTU (Maximum Transmission Unit)). They call flows "better than Poisson" (BTP) if the workload is larger when these flows are replaced by Poisson streams. This property has an important implication for studying overflow probabilities inside the network. Indeed, Massoulié shows in [28] via a sample-path Large Deviations (sp-LD) ordering that, if point processes satisfy the sp-LD principles with a finite rate function, the sp-LD dominance by the Poisson process is preserved in tandem FIFO queues. These desirable results thus inspire us to explore the BTP property of regulated flows.

Regulated flows do exemplify similar BTP properties. The authors of [19] observed that, as the number of sources increases with fixed total load, the mean delay of regulated traffic tends to converge to that of an $M/G/1$ queue fed with a marked Poisson process of parameters associated with the (σ, ρ, π) values; see Cao et al. [6] for results in a more general setting. Intuitively, when many of these essentially deterministic on-off processes are multiplexed together, the workload they generate is smoother than that from a Poisson process. However, Massoulié [28] also points out that deterministic processes do not have a smaller sp-LD rate function than a Poisson process for each sample path of their multiplexing. We thus need to define a weaker asymptotic ordering in which deterministic regulated flows are BTP.

Thus, in the second part of this paper, we first fix the total input load of regulated flows and study the asymptotic stochastic ordering between these fluid processes with a marked Poisson process, when the number of sources increases. With this scaling scheme, we find that the superposition of a large number of independent regulated flows is asymptotically smaller than a well-defined marked Poisson process. Furthermore, the workload generated by the superposed flow is also asymptotically dominated by an $M/G/1$ queue fed with the marked Poisson process. Besides this BTP property in terms of many sources with fixed load, we show that the large buffer decay rate and the many-sources asymptotic rate of regulated flows are larger than those of the $M/G/1$ queue. All these results establish the asymptotic BTP property of regulated flows.

The organization of this paper is as follows. In Section 2, we formulate our problem and present a preliminary upper bound for the freeze-out fraction when a single regulated flow accesses the queue. The statistical multiplexing of many streams are then considered in Section 3 where we obtain an upper bound via many-sources asymptotics. Subsequently, we study the asymptotic BTP property of regulated flows in Section 4. Finally, Section 5 presents our concluding remarks for the paper.

2. Preliminary results

We consider the scenario where M classes of independent regulated flows are multiplexed into a single FIFO queue with infinite buffer and server rate C . Each flow is leaky bucket regulated; a flow of class i is characterized by the $(\sigma_i, \rho_i, \pi_i)$ traffic descriptor, representing bucket size, average rate and peak rate, respectively. Let $A_i(s, t)$ denote the amount of data generated by a flow of class i in the time interval $(s, t]$. The quantity $A_i(s, t)$ satisfies for all $0 \leq s \leq t$

$$A_i(s, t) \leq \min(\pi_i(t - s), \sigma_i + \rho_i(t - s)).$$

There are Nn_i flows of class i and the total input rate $\rho = \sum_{i=1}^M Nn_i \rho_i < C$ such that the system is stable and a stationary regime exists. Also we need $\sum_{i=1}^M Nn_i \pi_i > C$ since otherwise the buffer would always be empty. Throughout this paper, we use a fluid model for input flows so that the processes $(A_i(0, t))$ are continuous and piecewise differentiable. We also assume that these processes are with ergodic and stationary increments.

To study the tail distribution of the workload in this queue, we first define the freeze-out fraction \mathcal{P} above a given threshold b as

$$\mathcal{P} \stackrel{\text{def}}{=} \lim_{t \rightarrow \infty} \frac{1}{A^N(t)} \int_0^t \mathbb{1}_{\{W^N(s) \geq b\}} A^N(ds),$$

where $W^N(t)$ is the amount of fluid in the buffer at time t and

$$A^N(t) = \sum_{i=1}^M \sum_{j=1}^{Nn_i} A_{i,j}(0, t) \tag{1}$$

is the total amount of fluid which arrives at the buffer in the time interval $[0, t]$. The quantity \mathcal{P} is the fraction of bits which enter the infinite buffer while the workload exceeds the threshold b . This quantity is a bound for the loss probability when the buffer capacity is finite and equal to b . With the knowledge of Palm theory [25] and the ergodicity of $(A^N(t))$, we know that $\mathcal{P} = C\mathbb{P}(W^N(0) \geq b)/\rho = \mathbb{P}(W^N(0) \geq b \mid W^N(0) > 0)$, where $\mathbb{P}(W^N(0) \geq b)$ is the probability that the workload in the queue exceeds b in the stationary regime. It is well known in the queueing literature that the evolution of the workload process $(W^N(t))$ is governed by the following differential equation:

$$dW^N(t) = -C\mathbb{1}_{\{W^N(t) > 0\}}dt + A^N(dt).$$

In the following, we begin with exhibiting an extremal traffic pattern for a single regulated source $(A(t))$, which maximizes the freeze-out fraction of the queue. This traffic profile later facilitates the study of the overflow probability of a queue when many flows are multiplexed together. The workload process in the queue is denoted by $(W(t))$ in the single input case.

Let us consider a busy period of the buffer and let τ denote the length of this busy period starting, say, at time 0. The freeze-out fraction $\mathcal{P}(\tau)$ over this busy period can be written as

$$\mathcal{P}(\tau) = \frac{1}{A(\tau)} \int_0^\tau \mathbb{1}_{\{W(s) \geq b\}} A(ds) = \frac{1}{C\tau} \int_0^\tau \mathbb{1}_{\{W(s) \geq b\}} A(ds).$$

We clearly have $\mathcal{P} = \mathbb{E}(\mathcal{P}(\tau))$. Our aim is to find the traffic profile which maximizes the quantity $\mathcal{P}(\tau)$ generated over the busy period considered. Instead of deriving the extremal traffic profile only for $\mathcal{P}(\tau)$, we can formulate this problem in a more general setting. For a given input process $(A(t))$, we define a functional $J(A)$ by

$$J(A) \stackrel{\text{def}}{=} \frac{1}{A(\tau)} \int_0^\tau f(W(s))A(ds),$$

where f is a monotonically increasing function. Our objective is to find a profile of $(A(t))$ such that the maximum of $J(A)$ is obtained over the given active period $[0, \tau]$, while the workload process $(W(t))$ is in the set \mathcal{Y} where

$$\mathcal{Y} = \left\{ W \in C_p^1[0, \tau] : W(0) = W(\tau) = 0 \right\},$$

with $C_p^1[0, \tau]$ denoting the set of functions which are continuous in $[0, \tau]$ and piecewise differentiable in $(0, \tau)$.

Lemma 1. The traffic pattern of the single source process $(A(t))$ which maximizes $J(A)$ is independent of the function f and is defined as follows:

- If $\tau \leq (\pi\sigma)/(C(\pi - \rho))$, the extremal traffic pattern is periodic and composed of a burst at the peak rate with duration $C\tau/\pi$, followed by a silence period with duration $(C - \pi)\tau/\pi$.
- If $\pi\sigma/(C(\pi - \rho)) \leq \tau \leq \tau_{\max} \stackrel{\text{def}}{=} \sigma/(C - \rho)$, the extremal traffic is periodic and composed of a burst at the peak rate π with length $\sigma/(\pi - \rho)$, followed by an activity period at rate ρ with length $C(\tau - \frac{\pi\sigma}{C(\pi-\rho)})/\rho$, and followed in turn by a silence period with length $(C - \rho)(\frac{\sigma}{C-\rho} - \tau)/\rho$.

Remark 1. The traffic pattern obtained in Lemma 1 maximizes a number of quantities related to queueing performances. When $f(x) = \mathbb{1}_{\{x \geq b\}}$, $J(A) = \mathcal{P}(\tau)$ and the freeze-out fraction is maximized. Taking $f(x) = x/C$, it can be computed as in [19] and $J(A)$ is associated with the average delay seen by arrivals. Other quantities such as the moment-generating functions of the workload can be studied by taking appropriate functions f .

Proof. First we define on \mathcal{Y} a partial order \leq by

$$W \leq V \quad \text{iff} \quad W(t) \leq V(t) \quad \text{for all } 0 \leq t \leq \tau.$$

Since f is an increasing function, it is readily seen that, if $(A(t))$ and $(B(t))$ are two arrival process giving rise to the workload processes $(W(t))$ and $(V(t))$, respectively, such that $W \leq V$, then $J(A) \leq J(B)$. Thus we need to find the extremal element $W^* \in \mathcal{Y}$ such that, for all $W \in \mathcal{Y}$, $W \leq W^*$. We can construct the following function $W^*(t)$:

- if $\tau \leq \pi\sigma/(C(\pi - \rho))$,

$$W^*(t) = \begin{cases} (\pi - C)t, & 0 \leq t \leq t'_1 \stackrel{\text{def}}{=} C\tau/\pi, \\ C(\tau - t), & t'_1 \leq t \leq \tau \end{cases}$$

- if $\tau \geq \pi\sigma/(C(\pi - \rho))$,

$$W^*(t) = \begin{cases} (\pi - C)t, & 0 \leq t \leq t_1 \stackrel{\text{def}}{=} \frac{\sigma}{\pi - \rho}, \\ \sigma + (\rho - C)t, & t_1 \leq t \leq t_2 \stackrel{\text{def}}{=} \frac{C\tau - \sigma}{\rho}, \\ \sigma + \rho t_2 - Ct, & t_2 \leq t \leq \tau. \end{cases}$$

Indeed, in the case $\tau \leq \pi\sigma/(C(\pi - \rho))$ (resp. $\tau \geq \pi\sigma/(C(\pi - \rho))$), owing to the (σ, ρ, π) constraint, $W(t) \leq W^*(t)$ for all $t \in [0, t'_1]$ (resp. $t \in [0, t_2]$). Now, assume that there exists some $t_0 \in [t'_1, \tau]$ (resp. $t_0 \in [t_2, \tau]$) such that $W(t_0) > W^*(t_0)$. Then,

$$-\int_{t_0}^{\tau} W'(s)ds > W^*(t_0) = C(\tau - t_0),$$

which implies that there exists a $t'_0 \in [t_0, \tau]$ such that $W'(t'_0) < -C$. This latter inequality is not possible since the drain rate from the queue cannot exceed C . As a consequence, for every $W \in \mathcal{Y}$, we have $W \leq W^*$.

Now, returning to the input process, when $\tau \leq \pi\sigma/(C(\pi - \rho))$, the input process, which maximizes the freeze-out fraction in the busy period with length τ , is the classical on-off process; during the on period the arrival rate is equal to the peak rate and the length of the on period is equal to $t'_1 = C\tau/\pi$.

In the case when $\tau \geq \pi\sigma/(C(\pi - \rho))$, the input process, which realizes the optimal trajectory $W^*(t)$ over a busy period, is composed of a burst at the peak rate π and with duration t_1 , followed by an activity period at rate ρ with length $(t_2 - t_1)$, and then by a silence period with length S given by $S = \tau - t_2$. Note that S is positive if and only if $\tau < \sigma/(C - \rho)$. The length of the busy period of a queue with an input process satisfying a (σ, ρ, π) -constraint is thus necessarily upper bounded by $\sigma/(C - \rho)$. \square

As a consequence of Lemma 1, the maximum freeze-out fraction $\mathcal{P}^*(\tau)$ is obtained from the extremal function W^* as

$$\mathcal{P}^*(\tau) = \begin{cases} 0 & \tau \leq \pi b/(C(\pi - C)) \\ 1 - \frac{\pi b}{C\tau(\pi - C)} & \pi b/(C(\pi - C)) \leq \tau \leq \tau_b^* \\ \mathcal{P}(\tau_b^*)\tau_b^*/\tau & \tau_b^* \leq \tau \leq \tau_{\max}, \end{cases}$$

and the maximum of $\mathcal{P}(\tau)$ is obtained at the critical length of the busy period $\tau_b^* = (C\sigma - \rho b)/(C(C - \rho))$ with the value $\mathcal{P}^*(\tau_b^*) = (\sigma - \frac{\pi-\rho}{\pi-C}b)/(\sigma - \frac{\rho}{C}b)$. Hence we obtain the upper bound for the freeze-out fraction over an arbitrary busy period of a queue fed with a single regulated input flow, which agrees with the conclusion in [22]. Since $\mathcal{P} = \mathbb{E}(\mathcal{P}(\tau))$, we obtain the following result.

Proposition 1 (Kesidis and Konstantopoulos [22]). *Under the assumptions $\pi > C > \rho$ and $(\pi - C)\sigma/(\pi - \rho) > b$, the freeze-out fraction \mathcal{P} in the single server queue fed with a (σ, ρ, π) -regulated fluid traffic source is upper bounded as*

$$\mathcal{P} \leq \frac{\sigma - \frac{\pi-\rho}{\pi-C}b}{\sigma - \frac{\rho}{C}b} \stackrel{\text{def}}{=} \mathcal{P}_{\max}. \tag{2}$$

The upper bound for the overflow probability in the stationary regime when the buffer capacity is finite and equal to b is then $\mathbb{P}(W(0) > b) \leq \rho\mathcal{P}_{\max}/C$.

As noted in [22,23], the worst-case traffic pattern is not unique. However, our worst-case profile also turns out to be natural in the subsequent analysis involving many sources, as discussed in Section 3.

To conclude this section, we establish the dominance of an $M/G/1$ queue over the queue fed with a (σ, ρ, π) -regulated flow. Let \tilde{W} denote the content in the stationary regime of a buffer drained at constant rate C and fed with batches with size σ arriving according to a Poisson process with intensity ρ/σ . The following proposition states that \tilde{W} is stochastically greater than W .

Proposition 2. *We have $W \leq_{st} \tilde{W}$, i.e., for all $b \geq 0$, $\mathbb{P}(W \geq b) \leq \mathbb{P}(\tilde{w} \geq b)$.*

Proof. From Proposition 1, we know that, for $b \leq \sigma$,

$$\mathbb{P}(W \geq b \mid W > 0) \leq \frac{\sigma - \frac{\pi-\rho}{\pi-C}b}{\sigma - \frac{\rho}{C}b} \leq \frac{\sigma - b}{\sigma - \frac{\rho}{C}b},$$

where the last inequality is obtained by letting $\pi \rightarrow \infty$. Now, by using a classical result by Erlang [30], we know that, over the interval $[j\sigma, (j + 1)\sigma]$ for $j \geq 0$,

$$\mathbb{P}(\tilde{W} \leq x) = \left(1 - \frac{\rho}{C}\right) \sum_{i=0}^j \frac{(i - x/\sigma)^i}{i!} \left(\frac{\rho}{C}\right)^i e^{-\frac{\rho}{C}(i-x/\sigma)}.$$

For $b \in [0, \sigma]$, it is easily checked that

$$\frac{\sigma - b}{\sigma - \frac{\rho}{C}b} \leq \frac{C}{\rho} \left(1 - \left(1 - \frac{\rho}{C}\right) \exp\left(\frac{\rho b}{C\sigma}\right)\right),$$

since $\frac{\rho b}{C\sigma} < 1$ and then, $\exp(\frac{\rho b}{C\sigma}) \leq \sigma/(\sigma - b\frac{\rho}{C})$. Hence, for all $b \geq 0$, $\mathbb{P}(W \geq b \mid w > 0) \leq \mathbb{P}(\tilde{W} \geq b \mid \tilde{W} > 0)$, and the result follows. \square

Proposition 2 implies that a single regulated flow is “better-than-Poisson”. We will extend this result in Section 4 to the better-than-Poisson property of regulated flows when many of them are multiplexed in a queue.

3. Many-sources asymptotics

We now consider the problem of estimating the freeze-out fraction when a large number of independent traffic streams are fed into a queue. When the transmission capacity C is large, and in particular, $C/\pi = O(N)$, we are in the regime of the many-sources asymptotics, which have been studied by many authors; see for instance [4,10,26]. In this section, we extend the results and formalism developed in Likhhanov and Mazumdar [26] to the continuous-time case by a discretization argument. Such extensions are similarly discussed in papers by Mandjes and Kim [27] and Guibert and Simonian [17], where they assume a local convex behavior of a rate function (see Eq. (3) below).

When a superposition of independent flows $(A^N(t))$ as defined by Eq. (1) enters a queue with server rate $C = Nc$, the stationary workload $(W^N(t))$ in the queue satisfies Reich’s formula: $W^N(0) = \sup_{t \geq 0} (A^N(-t, 0) - Nct)$. Let

$\phi_{i,t}(h) = \mathbb{E}(e^{hA_i(0,t)})$ denote the moment-generating function associated with a flow of class i where $h > 0$ and define the rate function associated with $A^N(0, t)$ by

$$I_t(a) = \sup_{h>0} \left(ah - \sum_{i=1}^M n_i \ln(\phi_{i,t}(h)) \right).$$

We now state the main result regarding the stationary tail distribution shown in [26] for the discrete-time case. This result can be extended to the continuous-time case by using the continuity of $I_t(a)$ with respect to a .

Proposition 3. Assume that there exists a unique $t_0 < \infty$ such that

$$I_{t_0}(ct_0 + b) = \min_{t \geq 0} I_t(ct + b) > 0. \tag{3}$$

Suppose that $\liminf_{t \rightarrow \infty} I_t(ct + b) / \log t > 0$ and that $I_t(x)$ is continuous in x . Then, as $N \rightarrow \infty$,

$$\mathbb{P}(W^N(0) > Nb) = \frac{e^{-NI_{t_0}(ct_0+b)}}{\tau \sqrt{2\pi\kappa^2 N}} \left(1 + O\left(\frac{1}{N}\right) \right), \tag{4}$$

where τ is the unique solution to the equation

$$\sum_{i=1}^M n_i \frac{\phi'_{i,t_0}(\tau)}{\phi_{i,t_0}(\tau)} = ct_0 + b$$

and

$$\kappa^2 = \left(\sum_{i=1}^M n_i \frac{\phi''_{i,t}(\tau)}{\phi_{i,t}(\tau)} - \left(\sum_{i=1}^M n_i \frac{\phi'_{i,t}(\tau)}{\phi_{i,t}(\tau)} \right)^2 \right).$$

It is important to note the existence of t_0 , which denotes the critical or most likely time-scale to overflow. Proposition 3 can be used to derive the worst-case traffic pattern for $(\sigma_j, \rho_j, \pi_j)$ -regulated flows in the case of many sources as follows.

Proposition 4. Let $r_i(t)$ be the rate function of the arrival process $(A_i(t))$ of a flow of class i , which is continuous and piecewise differentiable, such that $A_i(0, t) = \int_0^t r_i(s)ds \leq \min(\pi_i t, \rho_i t + \sigma_i)$ with $\sum_{i=1}^M n_i \rho_i < c$. The extremal sources which maximize the tail distribution of the queue are periodic on-off processes with random phases. In each period, the rate function $r_i(t)$ of source class i is given by

$$r_i(t) = \begin{cases} \pi_i, & 0 \leq t \leq t_i^1 \stackrel{\text{def}}{=} \sigma_i / (\pi_i - \rho_i), \\ \rho_i, & t_i^1 < t \leq t_i^1 + t_i^0 \stackrel{\text{def}}{=} t_0, \\ 0, & t_0 < t \leq t_0 + \sigma_i / \rho_i, \end{cases} \tag{5}$$

where t_0 is the most likely time-scale to overflow in Eq. (3) for the input process $(A^N(t))$ and t_i^0 denotes the duration of transmission at rate ρ_i which is determined by t_0 .

Proof. Note that $A^N(-t, 0) \stackrel{d}{=} A^N(0, t)$ due to stationarity. From the proof of the many-sources asymptotics in [26], it follows that, for some $\varepsilon > 0$,

$$\mathbb{P}(W^N(0) > Nb) = \mathbb{P}(A^N(t_0) > N(ct_0 + b)) \times (1 + O(e^{-\varepsilon N})).$$

In the following, we take t_0 to be fixed and investigate the contribution of the class i sources to the bound. Define $A_i^{N_i}(t) = \sum_{j=1}^{N_i} A_{i,j}(0, t)$ and $A^{N_i-}(t) = A^N(t) - A_i^{N_i}(t)$. In other words, we consider all inputs except the class i inputs. Then, denoting by $c_t = ct + b$, we have

$$\mathbb{P}(A^N(t_0) > Nct_0) = \int_0^\infty \mathbb{P}(A^{N_i-}(t_0) > Nct_0 - y) \times d\mathbb{P}(A_i^{N_i}(t_0) \leq y).$$

From the theorem of Bahadur–Rao [1], we obtain

$$\mathbb{P}(A^{N_i^-}(t_0) > Nct - y) = K(N)e^{\tau y} \left(1 + O\left(\frac{1}{N}\right) \right),$$

where $K(N) = e^{-NI_0(ct_0)} / (\tau\sqrt{2\pi\kappa^2N})$ as in Eq. (4) and does not depend on y , and τ is solution to the equation

$$\sum_{j=1, j \neq i}^M \frac{\phi'_{j,t_0}(\tau)}{\phi_{j,t_0}(\tau)} = ct_0 + b.$$

Hence, up to an error factor of $(1 + O(\frac{1}{N}))$, we obtain

$$\mathbb{P}(W^N(0) > Nb) = K(N) \int_0^\infty e^{\tau y} d\mathbb{P}(A_i^{N_i}(0, t_0) \leq y)$$

and then,

$$\mathbb{P}(W^N(0) > Nb) = K(N)(\mathbb{E}(e^{\tau A_i(0,t_0)}))^{Nn_i} = K(N)(\phi_{i,t_0}(\tau))^{Nn_i}. \tag{6}$$

From the ergodicity of the source, we have

$$\phi_{i,t_0}(\tau) = \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T e^{\tau A_i(s,s+t_0)} ds.$$

Hence it is clear that, in order to bound the quantity $\mathbb{P}\{W^N(0) > Nb\}$, it is sufficient to maximize the right-hand side of Eq. (6). We need determine the extremal traffic pattern which maximizes $\mathbb{E}(e^{\tau A_i(0,t_0)})$. With fixed t_0 and τ , since $A_i(0, t)$ is $(\sigma_i, \rho_i, \pi_i)$ -regulated with piecewise rate function, it is clear to see that Eq. (5) gives the extremal rate function which maximizes the quantity of information over a given period $[0, t_0]$ (see Remark 1). \square

With the above result, we can state the main result for the overflow probability.

Theorem 1. Consider a fluid queueing system with server rate $C = Nc$ and an infinite buffer which is accessed by $\sum_{i=1}^M Nn_i$ independent regulated sources. Assume that $\mathbb{E}(A_i(0, 1)) = \rho_i$ and $\sum_{i=1}^M n_i \rho_i < c$. Then, as $N \rightarrow \infty$

$$\mathbb{P}(W^N(0) > Nb) = \frac{e^{-NI_0(ct_0+b)}}{\tau\sqrt{2\pi\kappa^2N}} \left(1 + O\left(\frac{1}{N}\right) \right), \tag{7}$$

where the quantities $I_0(ct_0 + b)$, τ and κ are calculated as follows:

- Define

$$\phi_{i,t}(\tau) = \frac{1}{t + \frac{\sigma_i}{\rho_i}} \int_0^{t + \frac{\sigma_i}{\rho_i}} e^{\tau \int_u^{u+t} r_i(s) ds} du,$$

where $r_i(t)$ is the function defined in Proposition 4.

- Compute

$$I_0(ct_0 + b) = \inf_t \sup_h \left((ct + b)h - \sum_{i=1}^M n_i \ln(\phi_{i,t}(h)) \right). \tag{8}$$

- Compute τ as the solution to the equation

$$\sum_{i=1}^M n_i \frac{\phi'_{i,t_0}(\tau)}{\phi_{i,t_0}(\tau)} = ct_0 + b.$$

- Finally, compute

$$\kappa^2 = \sum_{i=1}^M n_i \frac{\phi''_{i,t_0}(\tau)}{\phi_{i,t_0}(\tau)} - \left(\sum_{i=1}^M n_i \frac{\phi'_{i,t_0}(\tau)}{\phi_{i,t_0}(\tau)} \right)^2.$$

Remark 2. The validity of our many-sources asymptotic holds when $\frac{Nc}{\pi} = O(N)$, i.e., when a large number of regulated flows are multiplexed in the queue and the buffer grows correspondingly. If the peak rates of some flows are comparable with the capacity C of the server, the Gaussian approximation in the theorem of Bahadur–Rao [1] does not hold. In this case, we can however use an upper bound for the overflow probability obtained via similar worst traffic profiles and Chernoff’s inequality as

$$\mathbb{P}(W^N(0) > b) \leq e^{-NI_0(Ct_0+b)}, \tag{9}$$

where the rate function I_0 can be computed as in Eq. (8); see also [32].

The principal difference with [32] is that, in our approach, we explicitly take into account the many-sources effect and determine the rate function for the source, which is extremal for the overflow asymptotics, rather than a priori first bounding the probability and then trying to make the bound small [8,32].

4. Better-than-Poisson asymptotics

Besides analyzing the tail distribution of a single queue with independent regulated input flows, as mentioned in the Introduction it is more interesting to characterize additional statistical properties of these flows such that their queuing performance inside a large network are numerically computable. Motivated by the better-than-Poisson conjecture studied in [3], in this section we compare the asymptotic workload distribution generated by regulated flows with that of an $M/G/1$ queue.

4.1. Many sources with fixed load

We begin with a scenario similar to that defined in Section 2, except that, as N increases, the average rate of the total aggregate ($A^N(t)$) defined by Eq. (1) is fixed and equal to $\rho = \sum_{i=1}^M Nn_i\rho_i < C$. Our goal is to identify the stochastic ordering between the workload ($W^N(t)$) generated by the superposition of regulated flows and the workload ($W^P(t)$) in the $M/G/1$ queue fed with a marked Poisson process ($A^P(t)$) with arrival rate $\lambda_p = \sum_{j=1}^M Nn_j\rho_j/\sigma_j$ and with marks B_p such that $\mathbb{P}(B_p = \sigma_i) = (Nn_i\rho_i)/(\sigma_i\lambda_p)$.

Theorem 2. Let each of the regulated flows ($A_{i,j}(t)$) assume the rate function given in (5) with

$$\lim_{\rho_i \rightarrow 0} \rho_i t_i^0 = 0, \tag{10}$$

then when the load ρ_0^i of each class $i = 1, \dots, M$, is fixed ($\rho_0^i = Nn_i\rho_i$ for class i), we have, for all $x \geq 0$,

$$\lim_{N \rightarrow \infty} \mathbb{P}(A^N(t) > x) \leq \mathbb{P}(A^P(t) > x) \tag{11}$$

and with regard to workload processes,

$$\lim_{N \rightarrow \infty} \mathbb{P}\{W^N(0) > x\} \leq \mathbb{P}\{W^P(0) > x\}. \tag{12}$$

Indeed, when the input load ρ is fixed, as N increases, the load of individual regulated sources decreases, which results in the bursts of a given source being more and more spread out. This intuitively explains why regulated flows tend to converge to a marked Poisson process from the viewpoint of the Poisson convergence theorem. Our main theorem rigorously shows that the multiplexing of independent regulated flows is asymptotically smaller than a marked Poisson process and the corresponding workload distribution is asymptotically dominated by an $M/G/1$ queue explicitly known. In the following, we first set up the ordering using the extremal profile constructed in Eq. (5) in Section 3.

The stochastic ordering between regulated flows ($A_i(t)$) and a marked Poisson process ($A_i^P(t)$) is set up via the construction of the following two marked point processes ($X_i(t)$) and ($X_i^P(t)$). As illustrated in Fig. 1, for each class i , the process ($X_i(t)$) has batch arrivals of size $B_i = \sigma_i + \rho_i t'_{0i}$ and inter-arrival time $\tau_i = t'_{0i} + \frac{\sigma_i}{\rho_i}$, where t'_{0i} (see Eq. (5))

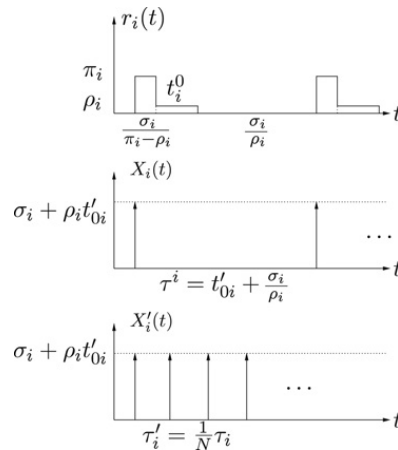


Fig. 1. Constructed processes for ordering.

is the length of the on period of $(A_i(t))$, with $t_{0i}' = \sigma_i / (\pi_i - \rho_i) + t_i^0$. In addition, let the process $(X_i'(t))$ have the same batch size B_i but inter-arrival time $\tau_i' = \frac{1}{N} \tau_i$. Hence, the process $(X_i(t))$ is a replica of the process $(X_i'(t))$ but the time-scale is dilated by N , i.e., $(X_i(t))$ and $(X_i'(\frac{t}{N}))$ have the same distribution.

From the construction of $(X_i(t))$, it is clear to see that the process $(X_i(t))$ dominates the regulated input $(A_i(t))$ for each sample path and thus for the corresponding workload distributions. Therefore, [Theorem 2](#) holds if

- (1) the process $(X^N(t))$, which is the superposition of Nn_i independent copies of the processes $(X_i(t))$ for $i = 1, \dots, M$, converges to $(A^P(t))$ in distribution,
- (2) if the workload (\hat{W}_i^N) generated by $(X^N(t))$ has a stationary distribution which converges to the workload (W_i^P) of the $M/G/1$ queue with input $(A^P(t))$.

In the next sections, we show these properties, which in turn proves [Theorem 2](#).

4.1.1. Convergence of $(X^N(t))$ to $(A^P(t))$

The convergence of $(X^N(t))$ to $(A^P(t))$ is stated in the following proposition.

Proposition 5. *When the load $\rho_0^i = Nn_i \rho_i$ of each class $i = 1, \dots, M$ is fixed, the process $(X^N(t))$ weakly converges to the process $(A^P(t))$ when N tends to infinity; this is denoted, for short, by $(X^N(t)) \Rightarrow (A^P(t))$.*

Proof. Let us first consider the point process $(\hat{X}_i(t))$ counting the jumps of process $(X_i(t))$. This process has mean intensity $\rho_i / (\sigma_i + \rho_i t_{0i}') = \rho_0^i / (Nn_i(\sigma_i + \rho_i t_{0i}'))$, where ρ_0^i is fixed. If we superpose Nn_i independent copies $(\hat{X}_{i,j}(t))$, $j = 1, \dots, Nn_i$, of the point process $(\hat{X}_i(t))$, then we are almost in the same situation as in [[13](#), Proposition 9.2.IV] and we can conclude that the superposed process $(\hat{X}_i^{Nn_i}(t))$ such that $\hat{X}_i^{Nn_i}(t) = \sum_{j=1}^{Nn_i} \hat{X}_{i,j}(t)$ weakly converges to the Poisson process with intensity ρ_0^i / σ_i . The present situation, however, is slightly different from [[13](#), Proposition 9.2.IV] since the average intensity of the point process $(\hat{X}_i(t))$ is not exactly equal to $\rho_0^i / (Nn_i \sigma_i)$, but there is a correcting term, which tends to 0 when N goes to infinity owing to Assumption (10). The proof of [[13](#), Proposition 9.2.IV] can then be readily adapted to give convergence of the superposed process to the corresponding Poisson process.

With regard to the marks, the jumps size of the process $(X_i(t))$ is equal to $B_i = \sigma_i + \rho_i t_{0i}' \rightarrow \sigma_i$ when $N \rightarrow \infty$, thanks to Assumption (10). Hence, when we consider the superposed process $(X_i^{Nn_i}(t))$ of Nn_i independent copies of the individual process $(X_i(t))$, we immediately deduce that $(X_i^{Nn_i}(t))$ weakly converges to the marked Poisson process with intensity ρ_0^i / σ_i and marks equal to σ_i .

When M classes of independent regulated flows $(X_i^{Nn_i})$, $i = 1, \dots, M$ are multiplexed together as $X^N(t) = \sum_{i=1}^M X_i^{Nn_i}(t)$, we obtain the Poisson convergence of the process $(X^N(t))$ by independence of the processes $(X_i^{Nn_i}(t))$

and the infinite divisibility of the Poisson process. In fact, for all $x > 0$, we have

$$\begin{aligned} \lim_{N \rightarrow \infty} \mathbb{P}\{X^N(t) \leq x\} &= \lim_{N \rightarrow \infty} \sum_{(k_1, \dots, k_M) \in \mathbb{N}^M} \mathbb{P}\left(X_i^{Ni} = k_i B_i, \sum_{i=1}^m k_i B_i \leq x\right) \\ &= \sum_{(k_1, \dots, k_M) \in \mathbb{N}^M: \sum k_i B_i \leq x} \prod_{i=1}^M \lim_{N \rightarrow \infty} \mathbb{P}(X_i^{Nn_i}(t) = k_i B_i) \\ &= \sum_{(k_1, \dots, k_M) \in \mathbb{N}^M: \sum k_i B_i \leq x} \prod_{i=1}^M \mathbb{P}(X_i^p(t) = k_i B_i) = \mathbb{P}(X^p(t) \leq x) \end{aligned}$$

by infinite divisibility of the Poisson process. \square

4.1.2. Convergence of $(\hat{W}^N(t))$ to $(W^p(t))$

Define the marked point process $X'(t) \stackrel{\text{def}}{=} \sum_{i=1}^M \sum_{j=1}^{n_i} X'_{i,j}(t)$, which corresponds to the superposition of n_i independent copies of the processes $(X'_i(t))$ for $i = 1, \dots, M$. If we consider N independent copies of the process $(X'(t))$, denoted by $(X^j(t))$ for $j = 1, \dots, N$, then we have $X^N(t) \stackrel{d}{=} \sum_{j=1}^N X^j(\frac{t}{N})$. In the following, we will first study the properties of $X'(t)$ which lead to the convergence result on the workload $(\hat{W}^N(t))$ of a queue fed with $(X^N(t))$.

We apply the same technique as in the paper by Cao and Ramanan [6]. For this purpose, we check preliminary properties for the processes under consideration. We first note that the process $(\hat{X}'(t))$ counting the jumps of the process $(X'(t))$ is a simple stationary point process with finite intensity

$$\lambda' = \sum_{i=1}^M \frac{N n_i \rho_i}{\sigma_i + \rho_i t_{0i}} \leq \lambda_p.$$

In addition, we clearly have, for all $\theta > 0$ and $t > 0$, $\mathbb{E}(e^{\theta X'(t)}) < \infty$.

Lemma 2. For all $x \geq C > \rho$, the function

$$A_1(x) \stackrel{\text{def}}{=} \liminf_{t \rightarrow 0} \sup_{\theta \in [0, \infty)} \left(x\theta - \frac{1}{t} \log \mathbb{E}(e^{\theta X'(t)})\right)$$

is such that $A_1(x) > 0$.

Proof. From the independence of the different processes $X'_{i,j}(t)$, we have, for all $\theta > 0$,

$$\frac{1}{t} \log \mathbb{E}(e^{\theta X'(t)}) = \frac{1}{t} \sum_{i=1}^M n_i \log \mathbb{E}(e^{B_i \theta \hat{X}'_i(t)}).$$

Since $\hat{X}'_i(t) = \lambda'_i t - \alpha_i + \mathbb{1}_{\{T_1^i \leq \alpha_i \tau'_i\}}$, where $\lambda'_i = 1/\tau'_i$, $\alpha_i = t/\tau'_i - \lfloor t/\tau'_i \rfloor$, and T_1^i is the time of the first point of $(\hat{X}'_i(t))$ after 0, we have

$$\begin{aligned} \frac{1}{t} \log \mathbb{E}(e^{\theta X'(t)}) &= \frac{1}{t} \sum_{i=1}^M n_i \log \mathbb{E}(e^{B_i \theta (\lambda'_i t + \mathbb{1}_{\{T_1^i \leq \alpha_i \tau'_i\}} - \alpha_i)}) \\ &= \frac{1}{t} \sum_{i=1}^M n_i (B_i (\lambda'_i t - \alpha_i) \theta + \log(e^{B_i \theta} \alpha_i + 1 - \alpha_i)) \\ &= \rho \theta - \frac{1}{t} \sum_{i=1}^M n_i (B_i \alpha_i \theta - \log(e^{B_i \theta} \alpha_i + 1 - \alpha_i)). \end{aligned}$$

Since $\alpha_i = t/\tau'_i - \lfloor t/\tau'_i \rfloor$, we have $\lim_{t \rightarrow 0} \frac{1}{t} \sum_{i=1}^M n_i B_i \alpha_i \theta / \tau'_i = \rho \theta$ and

$$\lim_{t \rightarrow 0} \sum_{i=1}^M n_i \frac{1}{t} \log(e^{B_i \theta} \alpha_i + 1 - \alpha_i) = \sum_{i=1}^M n_i \lambda'_i (e^{B_i \theta} - 1).$$

Therefore, for any $\theta > 0$,

$$\Lambda_1(x) \geq \liminf_{t \rightarrow 0} \left(x\theta - \frac{1}{t} \log \mathbb{E}(e^{\theta X'(t)}) \right) = x\theta - \sum_{i=1}^M n_i \lambda'_i (e^{B_i \theta} - 1) \stackrel{\text{def}}{=} f(\theta).$$

The derivative of the function $f(\theta)$ is equal to $f'(\theta) = x - \sum_{i=1}^M n_i \lambda'_i B_i e^{B_i \theta}$, and we note that $f(0) = 0$ and $f'(0) = x - \rho > 0$ for $x > \rho$. Hence, for sufficiently small θ , $f(\theta) > 0$, and we deduce that $\Lambda_1(x) > 0$. \square

Lemma 3. We have

$$\Lambda_2(C) \stackrel{\text{def}}{=} \liminf_{t \rightarrow \infty} \frac{t}{\log t} \sup_{\theta \geq 0} \left(C\theta - \frac{1}{t} \log \mathbb{E}(e^{\theta X'(t)}) \right) > 0.$$

Proof. Take $\theta_t = \beta \frac{\log t}{t}$ for some $\beta > 0$. We have

$$\begin{aligned} \Lambda_2(C) &\geq \liminf_{t \rightarrow \infty} \left(\frac{t}{\log t} C\theta_t - \frac{1}{\log t} \log \mathbb{E}(e^{\theta_t X'(t)}) \right) \\ &= C\beta - \limsup_{t \rightarrow \infty} \frac{1}{\log t} \log \mathbb{E}(e^{\theta_t X'(t)}). \end{aligned}$$

As in the proof of Lemma 2, we have

$$\begin{aligned} \log \mathbb{E}[e^{\theta_t X'(t)}] &= \sum_{i=1}^M n_i \log e^{B_i \theta_t (\lambda'_i t - \alpha_i)} \mathbb{E}(e^{B_i \theta_t \mathbb{1}_{\{T_i^i \leq \alpha_i \tau'_i\}}}) \\ &\leq \sum_{i=1}^M n_i \log e^{B_i \theta_t (\lambda'_i t + 1)} = \rho \beta \log t + \sum_{i=1}^M n_i B_i \theta_t. \end{aligned}$$

Hence,

$$\begin{aligned} \Lambda_2(C) &\geq C\beta - \limsup_{t \rightarrow \infty} \left(\frac{1}{\log t} \log \mathbb{E}(e^{\theta_t X'(t)}) \right) \\ &\geq C\beta - \limsup_{t \rightarrow \infty} \frac{1}{\log t} \left(\rho \beta \log t + \sum_{i=1}^M n_i B_i \beta \frac{\log t}{t} \right) = (C - \rho)\beta > 0. \quad \square \end{aligned}$$

With the above properties of process $(X'(t))$, we are now able to prove the main result of this section.

Proposition 6. For all $x \geq 0$, the stationary workloads $(\hat{W}^N(t))$ and $(W^P(t))$ are such that

$$\lim_{N \rightarrow \infty} \mathbb{P}(\hat{W}^N(0) > x) = \mathbb{P}(W^P(0) > x).$$

Proof. The stationary workload $\hat{W}^N(0) \stackrel{d}{=} \sup_{t \geq 0} (X^N(t) - Ct)$. Given $T \in [0, \infty)$, we can define the functional

$$F_T(f) \stackrel{\text{def}}{=} \sup_{t \in [0, T]} [f_t - Ct]$$

such that, when $f_t = X^N(t)$, we have, for any finite $T > 0$,

$$\mathbb{P}(\hat{W}^N(0) > x) \leq \mathbb{P}(F_T(X^N(t)) > x) + \mathbb{P} \left(\sup_{t \in [T, \infty)} (X^N(t) - Ct) > x \right).$$

Then, using the continuous mapping theorem, the \mathbb{P} -a.s. continuity of the projection operator $f \rightarrow f_T$ ([2, Theorems 5.1 and 15.1]), and Proposition 5, we know that, for all $T \geq 0$, $F_T(X^N(\cdot)) \Rightarrow F_T(A^P(\cdot))$ as $N \rightarrow \infty$. This implies that, for any $\varepsilon, x > 0$, there exists $N_0(\varepsilon) > 0$ such that, for all $N > N_0(\varepsilon)$,

$$\left| \mathbb{P} \left(\sup_{t \in [0, T]} (X^N(t) - Ct) > x \right) - \mathbb{P} \left(\sup_{t \in [0, T]} (A^P(t) - Ct) > x \right) \right| < \varepsilon/3. \tag{13}$$

Using Lemmas 2 and 3, and similar techniques as in [6, Section III], we can show that, for all $x > 0$ and $\varepsilon > 0$, there exist finite $T_X(\varepsilon)$ and $N_X(\varepsilon) > 0$ such that, for $T > T_X(\varepsilon)$ and $N > N_X(\varepsilon)$,

$$\mathbb{P} \left(\sup_{t \in [T, \infty)} (X^N(t) - Ct) > x \right) < \varepsilon/3. \tag{14}$$

From the infinite divisibility of the Poisson process, similar results hold when marked Poisson process $(A^P(t))$ is the input, i.e., for all $x > 0$ and $\varepsilon > 0$, there exist finite $T_p(\varepsilon)$ and $N_p(\varepsilon) > 0$ such that, for $T > T_p(\varepsilon)$ and $N > N_p(\varepsilon)$,

$$\mathbb{P} \left(\sup_{t \in [T, \infty)} (X_p(t) - Ct) > x \right) < \varepsilon/3. \tag{15}$$

We then obtain, for all $x > 0$ and $\varepsilon > 0$, by taking $T \geq \max(T_X(\varepsilon), T_p(\varepsilon))$ and $N > \max(N_0(\varepsilon), N_X(\varepsilon), N_p(\varepsilon))$

$$\begin{aligned} & |\mathbb{P}(\hat{W}^N(0) > x) - \mathbb{P}(W^P(0) > x)| \\ & \leq |\mathbb{P}(\hat{W}^N(0) > x) - \mathbb{P}(F_T(X^N(t)) > x)| + |\mathbb{P}(W^P(0) > x) - \mathbb{P}(F_T(A^P(t)) > x)| \\ & \quad + |\mathbb{P}(F_T(X^N(t)) > x) - \mathbb{P}(F_T(A^P(t)) > x)|, \end{aligned}$$

which implies that

$$\begin{aligned} & |\mathbb{P}(\hat{W}^N(0) > x) - \mathbb{P}(W^P(0) > x)| \\ & \leq \mathbb{P} \left(\sup_{t \in [T, \infty)} (X^N(t) - Ct) > x \right) + \mathbb{P} \left(\sup_{t \in [T, \infty)} (X_p(t) - Ct) > x \right) + \varepsilon/3 < \varepsilon, \end{aligned}$$

and the result follows. \square

Remark 3. We have so far established the asymptotic better-than-Poisson properties for traffic flows which have fixed $(\sigma_i, \rho_i, \pi_i)$ parameters and extremal profiles given in (5). In fact, from the construction of processes $(X_i(t))$ and $(X'_i(t))$, it is clear to see that these properties can be extended to multiplexing regulated flows $(A_i(t))$ with general deterministic on-off profiles. Changing the peak rate constraint and/or on period profiles does not affect the convergence results. All that matters is that each flow $(A_i(t))$ has sustainable average rate $\rho_i = \lim_{t \rightarrow \infty} A_i(t)/t$ and its burstiness is confined as $\sup_{t > 0} (A_i(s, t+s) - \rho_i t) \leq \sigma_i$. The dominant marked Poisson process can be constructed with these (σ_i, ρ_i) parameters, as given in Theorem 2.

The dominance performance of the $M/G/1$ queue is illustrated in Fig. 2. We also compare against simulation the other performance upper bounds obtained via the many-sources asymptotic approach in Section 3 and the Bernoulli approach [8] which bounds the moment-generating function of regulated flows by a Hoeffding inequality. In the simulations, we fixed the capacity $C = 1$, the threshold $b = 0.3$ and the offered load $\rho = 0.7$ as the number of sources increases. We simulated the overflow probability of the queue when multiple homogeneous independent regulated flows are multiplexed together. Note that the peak rate $(\pi = 0.5, 2)$ of each flow is comparable with the capacity in our scenarios, so we chose Eq. (9) in Section 3 as the many-sources upper bound. The overflow probability of the corresponding $M/G/1$ queue can be computed via the inversion of its Laplace transform which has a well-known explicit form [24].

Fig. 2 clearly shows that our $M/G/1$ bound is comparable to the bounds from many-sources asymptotics. It becomes more accurate as the peak rate of each flow or the number of sources grows, i.e., when the regulated flows converge to the constructed marked point processes $(X^N(t))$. Moreover, although it is only the asymptotic ordering that we can show between regulated flows and Poisson, the numerical evaluation indicates that the workload increases as N increases, and the better-than-Poisson property holds for any N .

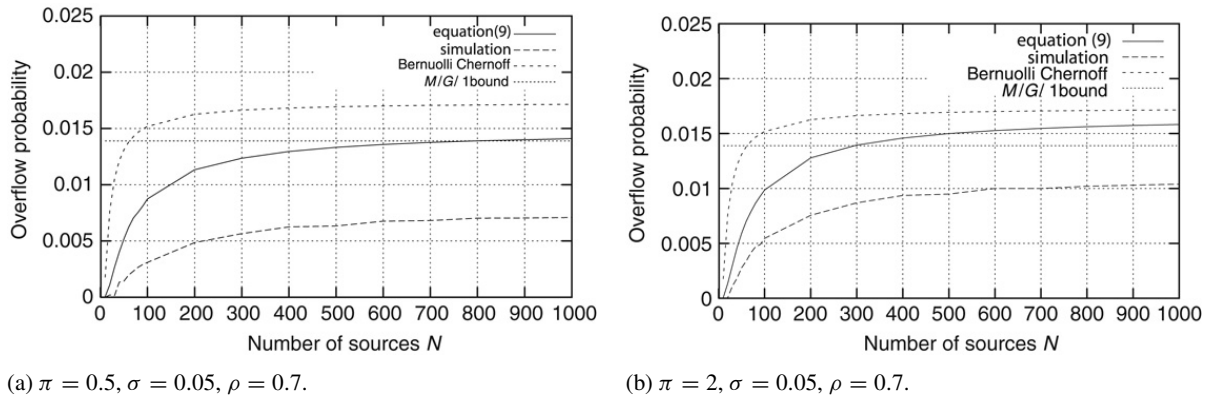


Fig. 2. Comparison of many sources bound, Bernoulli bound, and $M/G/1$ bound vs. simulation for independent homogeneous regulated sources.

4.2. Discussions on other asymptotics

Besides the better-than-Poisson asymptotics for many regulated flows with their total load fixed, we have also studied in [18] that regulated flows are better-than-Poisson in the large buffer asymptotics. We specifically have the following result.

Proposition 7. *Let the server rate be $C = Nc$ and the buffer threshold be $B = Nb$. Then the workload $W^N(0)$ in the queue with regulated inputs $(A^N(t))$ defined by Eq. (1) decays at a faster rate than the workload $W^P(0)$ of the $M/G/1$ queue, i.e.,*

$$\lim_{b \rightarrow \infty} \frac{1}{Nb} \log \mathbb{P}(W^N \geq Nb) \leq -\xi \tag{16}$$

where ξ is the Kingman's exponent of the $M/G/1$ queue, defined by $-\xi = \lim_{x \rightarrow \infty} \frac{1}{x} \log \mathbb{P}(W^P(0) \geq x)$.

Proof. We consider a partition $\{c_i\}, i = 1, \dots, M$, of the server capacity $C = Nc$ such that $\sum_{i=1}^M n_i c_i = c$ and $c_i > \rho_i$ for $i = 1, \dots, M$. From Reich's formula, we have

$$W^N(0) \stackrel{d}{=} \sup_{t \geq 0} \sum_{i=1}^M \sum_{j=1}^{Nn_i} (A_{i,j}(0, t) - c_i t) \leq \sum_{i=1}^M \sum_{j=1}^{Nn_i} W_{i,j}(0),$$

where $W_{i,j}(0)$ is the content in the stationary regime of a buffer drained at constant rate c_i and fed with a source of class i , which is $(\sigma_i, \rho_i, \pi_i)$ -regulated. Hence, for all $b \geq 0$, we have, from Proposition 2,

$$\mathbb{P}(W \geq Nb) \leq \mathbb{P}\left(\sum_{i=1}^M \sum_{j=1}^{Nn_i} W_{i,j} \geq Nb\right) \leq \mathbb{P}\left(\sum_{i=1}^M \sum_{j=1}^{Nn_i} \tilde{W}_{i,j}(0) \geq Nb\right),$$

where $\tilde{W}_{i,j}(0)$ is the stationary content in a fluid buffer drained at constant rate c_i and fed with batches of fluid with size σ_i arriving according to a Poisson process with rate ρ_i/σ_i . Its Laplace transform $\mathbb{E}(e^{-s\tilde{W}_{i,j}(0)})$ is given by

$$\mathbb{E}(e^{-s\tilde{W}_{i,j}}) = \frac{(c_i - \rho_i)s}{c_i s + \frac{\rho_i}{\sigma_i} e^{-\sigma_i s} - \frac{\rho_i}{\sigma_i}}.$$

Let $\tilde{s}_i(c_i)$ denote the module of the pole with the smallest module of the above Laplace transform. (The positive real number $\tilde{s}_i(c_i)$ is called the Kingman's exponent of the corresponding $M/G/1$ queue.) We have $\tilde{s}_i = \frac{1}{\sigma_i} \eta(\frac{c_i}{\rho_i})$, where for $x = c_i/\rho_i > 1$, $\eta(x)$ is the non-zero root with the smallest module of the equation

$$-x\eta + e^\eta - 1 = 0.$$

Note that $\eta(x)$ is an increasing positive function of x and it satisfies $\eta > 1 - 1/x$ and $\eta e^\eta > 2(e^\eta - x)$.

Now with independent inputs, the different random variables $\tilde{W}_{i,j}(0)$ are independent and the Laplace transform

$$\mathbb{E} \left(e^{-s \sum_{i=1}^M \sum_{j=1}^{Nn_i} \tilde{W}_{i,j}(0)} \right) = \prod_{i=1}^M (\mathbb{E}(e^{-s \tilde{W}_i}))^{Nn_i}, \tag{17}$$

where for fixed i , \tilde{W}_i is distributed as the random variables $\tilde{W}_{i,j}(0)$. The tail of the probability distribution function of the random variable $\sum_{i=1}^M \sum_{j=1}^{Nn_i} \tilde{W}_{i,j}(0)$ is governed by the root with the smallest module of the above Laplace transform; this root is $-\inf_i \{\tilde{s}_i(c_i)\}$. This property holds for any partition $\{c_i\}$ of the server capacity such that $c_i > \rho_i$ for all $i = 1, \dots, M$.

To determine the exact tail behavior of the probability distribution of the random variable W , we are led to determine the maximum value of $\inf\{\tilde{s}_i(c_i)\}$ over all the partitions of the server capacity C . Since $s(c_i) = \frac{1}{\sigma_i} \eta(\frac{c_i}{\rho_i})$ and $\eta(x)$ is an increasing function of x , the optimal values c_i^* are such that all the values $\frac{1}{\sigma_i} \eta(\frac{c_i^*}{\rho_i})$ for $i = 1, \dots, M$ are equal to some constant, say s^* . Indeed, if all the $s_i(c_i)$ were not equal, it would always be possible to increase the maximum value of the minimum by decreasing the largest value. Hence, s^* satisfies $-\frac{c_i^*}{\rho_i} \sigma_i s^* = e^{-\sigma_i s^*} - 1$, for $i = 1, \dots, M$. With $C = \sum_{i=1}^M Nn_i c_i^*$, s^* turns out to be the solution with the smallest module to the equation

$$Cs^* + \sum_{i=1}^M Nn_i \frac{\rho_i}{\sigma_i} (e^{-s^* \sigma_i} - 1) = 0.$$

Note that the Laplace transform for the workload $W^P(0)$ of the $M/G/1$ queue is

$$\mathbb{E}(e^{-s \tilde{W}}) = \frac{(C - \rho)s}{Cs + \sum_{i=1}^M Nn_i \frac{\rho_i}{\sigma_i} (e^{-\sigma_i s} - 1)}.$$

Therefore, $s^* = \xi$ is the Kingman's exponent of this $M/G/1$ queue with Poisson input with intensity $\lambda = \sum_{i=1}^M Nn_i \rho_i / \sigma_i$.

Going back to Eq. (17), since all the parameters $\tilde{s}_i(c_i^*)$ are equal to ξ , the point $-\xi$ is a pole with order $L \stackrel{\text{def}}{=} N \sum_{i=1}^M n_i$ for the Laplace transform (17). The Laplace transform $\mathbb{E}(e^{-s \tilde{w}_i})$ can specifically be written as

$$\mathbb{E}(e^{-s \tilde{W}_i}) = a_i \left(\frac{\xi}{s + \xi} - f_i(s) \right),$$

with $a_i = (c_i^* - \rho_i) / (\rho_i e^{\sigma_i \xi} - c_i^*)$ and

$$f_i(s) = \frac{c_i^* - \rho_i e^{\sigma_i \xi} + \xi \frac{\rho_i}{\sigma_i} \sum_{n=2}^{\infty} \frac{(-\sigma_i)^n}{n!} (s + \xi)^{n-2} e^{\sigma_i \xi}}{c_i^* - \rho_i e^{\sigma_i \xi} + \frac{\rho_i}{\sigma_i} \sum_{n=2}^{\infty} \frac{(-\sigma_i)^n}{n!} (s + \xi)^{n-1} e^{\sigma_i \xi}}.$$

It then follows that the Laplace transform (17) can be written as

$$\mathbb{E} \left(e^{-s \sum_{i=1}^M \sum_{j=1}^{Nn_i} \tilde{W}_{i,j}(0)} \right) = \left(\prod_{i=1}^M a_i^{Nn_i} \right) \left(\sum_{j=1}^L \kappa_j \left(\frac{\xi}{s + \xi} \right)^j + g(s) \right),$$

where function g has poles with modules greater than $-\xi$ and κ_j is the coefficient of Y^j in the product $\prod_{j=1}^L (Y + b_j)$, with $b_j = \xi \rho_i \sigma_i e^{\sigma_i \xi} / 2(\rho_i e^{\sigma_i \xi} - c_i^*) - 1$. Note that $b_j > 0$ and $a_i > 0$ from the properties of η .

By using [20, Theorem 10.7], it follows that

$$\mathbb{P} \left(\sum_{i=1}^M \sum_{j=1}^{Nn_i} \tilde{W}_{i,j}(0) \geq x \right) \sim \left(\prod_{i=1}^M a_i^{Nn_i} \right) \left(\sum_{j=1}^L \frac{\kappa_j}{(j-1)!} \Gamma(j, x\xi) \right),$$

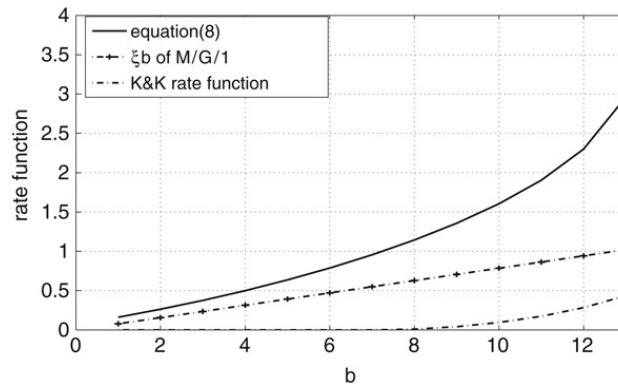


Fig. 3. Comparison of many-sources rate functions for homogeneous sources (i.e., $M = 1$) with parameters $\sigma = 16$, $\rho = 1$, $C = N \times 2$, and $\pi = 9$ —K&K rate function is the rate function obtained by Kesidis and Konstantopoulos [22].

where we have used the incomplete Gamma function $\Gamma(a, x) \stackrel{\text{def}}{=} \int_x^\infty t^{a-1} e^{-t} dt$. By using the asymptotic equivalent $\Gamma(a, z) \sim z^{a-1} e^{-z}$ for large $z > 0$, we obtain

$$\mathbb{P} \left(\sum_{i=1}^M \sum_{j=1}^{Nn_i} \tilde{W}_{i,j}(0) \geq x \right) \sim \frac{\kappa_L}{(L-1)!} \left(\prod_{i=1}^M a_i^{Nn_i} \right) (Nb\xi)^{L-1} e^{-Nb\xi}. \tag{18}$$

Take logarithms and divide by Nb ; Eq. (16) then follows by letting $b \rightarrow \infty$. \square

By presenting Proposition 7 in a similar scaling setting as the many-sources asymptotics in Section 3, we make it easier to study the relationship between the Kingman’s exponent ξ and the many-sources asymptotic rate $I_{t_0}(ct + b)$. In fact, the authors in [4,10] have studied the connections between the many-sources asymptotic rate $I_{t_0}(ct + b)$ and the large buffer asymptotic decay rate $H(c)$, defined as in [4]:

$$H(c) \stackrel{\text{def}}{=} - \lim_{b \rightarrow \infty} \frac{1}{Nb} \log \mathbb{P}\{W^N(0) \geq Nb\}.$$

They show that for the multiplexing of general independent traffic flows, as the buffer level b becomes large, these two asymptotics in our setting are equivalent in the limit sense, i.e.,

$$\lim_{b \rightarrow \infty} \frac{1}{b} I_{t_0}(ct + b) = H(c).$$

In addition, the authors of [4,10] also indicated that Poisson arrivals do not gain scale economies, i.e., $\frac{1}{b} I_{t_0}^P(ct + b) = \xi$. Therefore we can conclude that the many-sources asymptotic rate of regulated flows $I_{t_0}^r(ct + b)$ also dominates the Kingman’s exponent of the $M/G/1$ queue when the buffer size goes large, i.e.,

$$\lim_{b \rightarrow \infty} \frac{1}{b} I_{t_0}^r(ct + b) \geq \xi. \tag{19}$$

In fact, numerical computations show that $I_{t_0}^r(ct + b) \geq I_{t_0}^P(ct + b) = \xi b$ for any buffer level b , as seen in Fig. 3. However, a rigorous proof of this ordering property requires explicit and accurate characterization of the moment-generating function of the regulated flows, which demands further investigations.

All these better-than-Poisson discussions indicate that an appropriate $M/G/1$ queue can well estimate the workload distribution of regulated flows in a large network system. The performance bound given by this $M/G/1$ queue is more accurate when the peak rate of each regulated flow is comparable to the service capacity, i.e., when the many-sources asymptotic bound in Theorem 1 does not hold.

5. Conclusion

In this paper we have studied the tail distribution of the buffer content when independent heterogeneous $(\sigma_i, \rho_i, \pi_i)$ -regulated traffic streams are multiplexed in a FIFO manner. We obtained two types of upper bounds via the many-sources asymptotic and the asymptotic dominance property of an $M/G/1$ queue, respectively. Numerical results

indicate that these bounds function quite well and that substantial multiplexing gains are achievable although the input flows are simply characterized with three leaky bucket parameters.

Our many-sources asymptotic bound is obtained via identifying the worst-case traffic profiles of regulated flows. It is much tighter than the approximations given in the literature. Moreover, since the bound is achieved via the extremal sources for a (σ, ρ, π) envelope, it is unlikely that we can do much better without more information about the moment-generating function of the sources. Meanwhile, the bound via the “asymptotic better-than-Poisson” property provides insights in understanding the workload distribution from another perspective. It implies that the workload generated by fluid-regulated flows can be stochastically dominated by some $M/G/1$ queue. Both of these results complement those reported in [19] in that we can compute bounds for both the mean buffer occupancy as well as the asymptotic. It is also interesting to note that the extremal sources for both results have the same behavior. We can utilize the mean delay results for network design when the traffic is best effort while the results reported in this paper are better for tight quality of service constraints [16].

As claimed in [3], we expect to analyze the performance of regulated flows inside a network via constructing proper $M/G/1$ queues for performance upper bounds. Indeed, in [34], the authors found that, in a large network, each internal flow with initially fixed burst size σ_i can be regulated by the same pair of (σ_i, ρ_i) parameters. Thus each aggregate flow consisting of multiple individual flows from the same previous queues is regulated by $(\sum \sigma_i, \sum \rho_i)$. In addition, since flows from different previous routes can be regarded independent, these flows to an internal queue can thus be asymptotically dominated by a marked Poisson process constructed according to the (σ_i, ρ_i) parameters. This is a simplest way to construct an $M/G/1$ queue to dominate the workload distribution of internal queues with regulated inputs. However, this type of $M/G/1$ queue gives conservative estimates since it ignores the possible multiplexing gains in aggregate flows. One extension of this work is to find proper burstiness characteristics of aggregate flows inside the network, which should also be practical to compute for engineering purposes.

References

- [1] R.R. Bahadur, R. Ranga Rao, On deviations of the sample mean, *Ann. Math. Statist.* 31 (1960) 1015–1027.
- [2] P. Billingsley, *Convergence of Probability Measures*, John Wiley, New York, 1968.
- [3] T. Bonald, A. Proutière, J. Roberts, Statistical performance guarantees for streaming flows using expedited forwarding, in: *Proceedings of the IEEE INFOCOM 2001*, Alaska, USA, April 2001.
- [4] D.D. Botvich, N.G. Duffield, Large deviations, the shape of the loss curve, and economies of scale in large multiplexers, *Queueing Syst.* 20 (3–4) (1995) 293–320.
- [5] J.-Y. Le Boudec, Application of network calculus to guaranteed service networks, *IEEE Trans. Inform. Theory* 44 (3) (1998) 1087–1096.
- [6] J. Cao, K. Ramanan, A Poisson limit for buffer overflow probabilities, in: *Proceedings of INFOCOM 2002*, vol. 2, June 2002, pp. 994–1003.
- [7] C.-S. Chang, *Performance Guarantees in Communication Networks*, Springer-Verlag, London, 2000.
- [8] C.-S. Chang, W. Song, Y.-M. Chiu, On the performance of multiplexing independent regulated inputs, in: *Proceedings of Sigmetrics 2001*, Boston, MA, USA, May 2001.
- [9] F. Ciucu, A. Burchard, J. Liebeherr, A network service curve approach for the stochastic analysis of networks, in: *Proc. ACM Sigmetrics*, Banff, Alberta, Canada, June 6–10, 2005.
- [10] C. Courcoubetis, R.R. Weber, Buffer overflow asymptotics for a buffer handling many traffic sources, *J. Appl. Probab.* 33 (3) (1996) 886–903.
- [11] R.L. Cruz, A calculus for network delay. I. Network elements in isolation, *IEEE Trans. Inform. Theory* 37 (1) (1991) 114–131.
- [12] R.L. Cruz, A calculus for network delay. II. Network analysis, *IEEE Trans. Inform. Theory* 37 (1) (1991) 132–141.
- [13] D.J. Daley, D. Vere-Jones, *An Introduction to the Theory of Point Processes*, in: *Springer Series in Statistics*, Springer Verlag, 1988.
- [14] D. Eun, N.B. Shroff, Simplification of network analysis in large-bandwidth systems, in: *Proc. of IEEE INFOCOM 2003*, San Francisco, CA, March 2003.
- [15] V. Firoiu, J.-Y. Le Boudec, D. Towsley, Z.L. Zhang, Theories and models for internet quality of service, *Proceedings of the IEEE* 90 (9) (2002) 1565–1591.
- [16] A. Girard, C. Rosenberg, H. Cho, Optimal performance partitioning for networks with envelope-regulated traffic, in: *Proc ITC Specialist Semina 15*, Würzburg, Germany, July 2002.
- [17] J. Guibert, A. Simonian, Large deviations approximations for fluid queues fed by a large number of on/off sources, *IEEE Sel. Areas Commun.* 13 (6) (1995) 1017–1026.
- [18] F. Guillemin, N. Likhanov, R. Mazumdar, C. Rosenberg, Y. Ying, Buffer overflow bounds for multiplexed regulated traffic streams, in: *Proc. ITC 18*, Elsevier Science, Berlin, 2003.
- [19] F.M. Guillemin, N. Likhanov, R.R. Mazumdar, C.P. Rosenberg, Extremal traffic and bounds for the mean delay of multiplexed regulated traffic streams, in: *Proceedings of the IEEE INFOCOM 2002*, New York, USA, June 2002.
- [20] P. Henrici, *Applied and Computational Complex Analysis*, vol. 2, Wiley, New York, 1977.
- [21] W. Hoeffding, Probability inequalities for sums of bounded random variables, *Amer. Statist. Assoc. J.* (1963) 13–30.
- [22] G. Kesidis, T. Konstantopoulos, Extremal traffic and worst-case performance for queues with shaped arrivals, in: *Analysis of Communication Networks: Call Centres, Traffic and Performance*, Toronto, ON, 1998, Amer. Math. Soc., Providence, RI, 2000, pp. 159–178.

- [23] G. Kesidis, T. Konstantopoulos, Worst case performance of a buffer with independent shaped arrival processes, *IEEE Commun. Lett.* 4 (1) (2000) 26–28.
- [24] L. Kleinrock, *Queueing Systems Vol. 1: Theory*, Wiley Interscience, 1975.
- [25] T. Konstantopoulos, G. Last, On the dynamics and performance of stochastic fluid systems, *J. Appl. Probab.* 37 (3) (2000) 652–667.
- [26] N. Likhhanov, R.R. Mazumdar, Cell loss asymptotics for buffers fed with a large number of independent stationary sources, *J. Appl. Probab.* 36 (1) (1999) 86–96.
- [27] M. Mandjes, J.H. Kim, Large deviations for small buffers: an insensitivity result, *Queueing Syst.* 37 (4) (2001) 349–362.
- [28] L. Massoulié, Large deviations ordering of point processes in some queueing networks, *Queueing Syst.* 28 (4) (1998) 317–335.
- [29] L. Massoulié, A. Busson, Stochastic majorization of aggregates of leaky-bucket constrained traffic streams. Preprint, Microsoft Research, Cambridge, 2000.
- [30] R. Syski, *Introduction to Congestion Theory in Telephone Systems*, Oliver and Boyd, 1967.
- [31] M. Vojnovic, J.-Y. Le Boudec, Bounds for independent regulated inputs multiplexed in a service curve network element, *IEEE Trans. Commun.* 51 (5) (2003) 735–740.
- [32] M. Vojnovic, J.-Y. Le Boudec, Stochastic analysis of some expedited forwarding networks, in: *Proceedings of the IEEE INFOCOM 2002*, New York, USA, June 2002.
- [33] D. Wischik, The output of a switch, or, effective bandwidths for networks, *Queueing Syst.* 32 (1999) 383–396.
- [34] Y. Ying, R. Mazumdar, C. Rosenberg, F. Guillemin, Burstiness behavior of regulated flows inside networks, in: *Proceedings of Networking 2005*, Waterloo, Canada, IFIP, Springer, May 2005, pp. 918–929.



Yu Ying received her M.S.E.E. and B.S.E.E. degrees from Tsinghua University, Beijing, China in 2000 and 1998, respectively, and her Ph.D. in Electrical and Computer Engineering from Purdue University, West Lafayette, IN, USA in 2006. She is currently a Member of Technical Staff at Juniper Networks, Sunnyvale, CA, USA.

Her research interests are in QoS performance modeling and analysis of broadband networks and queueing theory.



Fabrice Guillemin received degrees from Ecole Polytechnique, Paris, in 1987, and from Ecole Nationale Supérieure des Télécommunications, Paris, in 1989. He defended his Ph.D. thesis in 1992 at the University of Rennes 1 and his habilitation thesis in 1999 at the University of Paris 6. Since 1989, he has been with the research and development center of France Telecom in Lannion, France. His research interests are in the performance evaluation of packet telecommunication networks (especially metrology of the Internet) and queueing theory, together with applied probability.



Ravi Mazumdar was born in Bangalore, India. He obtained his B.Tech. in Electrical Engineering from the Indian Institute of Technology, Bombay, India, in 1977, his M.Sc. DIC in Control Systems from Imperial College, London, UK, in 1978 and his Ph.D. in Systems Science from the University of California, Los Angeles, USA, in 1983.

He is currently a University Research Chair Professor of Electrical and Computer Engineering at the University of Waterloo, Waterloo, Canada, and an Adjunct Professor of ECE at Purdue University. He has served on the faculties of Columbia University (NY, USA), INRS-Telecommunications (Montreal, Canada), University of Essex (Colchester, UK), and most recently at Purdue University (West Lafayette, USA). He has held visiting positions and sabbatical leaves at UCLA, the University of Twente (Netherlands), the Indian Institute of Science (Bangalore), and the Ecole Nationale Supérieure des Télécommunications (Paris). He is an Associate Editor of the *IEEE/ACM Transactions on Networking*.

He is a Fellow of the IEEE and the Royal Statistical Society. He is a member of the working groups WG6.3 and 7.1 of the IFIP and a member of SIAM and the IMS. He won the Best Paper Award at INFOCOM 2006 and was co-author a paper that was runner-up for the Best Paper at INFOCOM 1998.

His research interests are in wireless and wireline networks, applied probability, queueing theory, and stochastic analysis with applications to traffic engineering, stochastic filtering theory, and mathematical finance.



Catherine Rosenberg was educated in France (Ecole Nationale Supérieure des Télécommunications de Bretagne, Diplôme d'Ingénieur in EE in 1983, and University of Paris, Orsay, Doctorat en Sciences in CS in 1986) and in the USA (UCLA, MS in CS in 1984). Dr. Rosenberg has worked in several countries including USA, UK, Canada, France and India. In particular, she has worked for Nortel Networks in the UK, AT&T Bell Laboratories in the USA, Alcatel in France and has taught at Purdue University (USA) and Ecole Polytechnique of Montreal (Canada).

Dr. Rosenberg is currently a University Research Chair Professor in the Department of Electrical and Computer Engineering at the University of Waterloo, Canada, where she also served as Department Chair from 2004–07. She is a Senior Member of the IEEE and Associate Editor of the *IEEE Trans. on Mobile Computing* and has been past editor of *Telecommunication Systems*. She was appointed to a three-year term on the Scientific Advisory Board of France Telecom in 2006.

Her research interests are broadly in networking with currently an emphasis in wireless networking and in traffic engineering (Quality of Service, Network Design, and Routing). She has authored over 70 papers and has been awarded seven patents in the USA.

Agencies and industries that have supported her research include USA NSF (National Science Foundation) NSERC (Natural Sciences and Engineering Research Council of Canada), FCAR (The Quebec counterpart of NSERC), CRC (Canadian Ministry of Communications), EEC (European Commission) while at Nortel Networks, ESA (European Space Agency) while at Nortel Networks, France-Telecom, CISCO, Bell Canada, and Nortel Networks.